

# **Investigative Data Mining Toolkit: A Software Prototype for Visualizing, Analyzing and Destabilizing Terrorist Networks**

**Nasrullah Memon and Henrik Legind Larsen**

Software Intelligence Security Research Center  
Department of Computer Science and Engineering  
Aalborg University, Niels Bohrs Vej 8  
6700 Esbjerg, Denmark  
[nasrullah@cs.aau.dk](mailto:nasrullah@cs.aau.dk), [legind@cs.aau.dk](mailto:legind@cs.aau.dk)

## **ABSTRACT**

Knowledge about the structure and organization of terrorist networks is important for both terrorism investigation and the development of effective strategies to prevent terrorists' attacks. However, except for network visualization, terrorist network analysis remains primarily a manual process. Existing tools do not provide advanced structural analysis techniques that allow extraction of network knowledge from large volumes of criminal-justice data. It is a well known fact that terrorist activities consist of dispersed organizations (like non-hierarchical organizations), small groups, and individuals who communicate, coordinate and conduct their campaign in a network-like manner. There is a pressing need to automatically collect data of terrorist networks, analyze such networks to find hidden relations and groups, prune datasets to locate regions of interest, find key players, characterize the structure, trace point of vulnerability, and detect efficiency of the network. To meet this challenge, we designed and developed a knowledgebase for storing and manipulating data collected from various authenticated websites. This paper presents framework of investigative data mining toolkit, our recently introduced techniques and algorithms (which are implemented in the investigative data mining toolkit) could be useful for law enforcement agencies that need to analyze terrorist networks and prioritize their targets. Applying recently introduced algorithms for constructing hidden hierarchy of non-hierarchical terrorist networks, we present case studies of the terrorist attacks that occurred in past, in order to construct command structure of the networks.

## **KEYWORDS**

Visualization, Analysis and Destabilizing Terrorist Networks, Position Role Index, Dependence Centrality, Hierarchy of Non-Hierarchical Networks

## **1.0 INTRODUCTION**

The threats facing society today require new methods for modeling and analysis. In fact, civil security decision makers, analysts and field operators fighting terrorism and organized crime across the European Union all need front-line integrated technologies to support their cooperative work. Our opponents are no longer organized in hierarchical structures, but instead consist of individuals and groups that are loosely organized in "dark networks". Instead of large-scale military attacks, they stage attacks or set bombs against unprotected civilians, or seek to influence crowds of legitimate demonstrators so that critical riot situations occur.

In order to construct decision support systems that take account of these new factors, new, more powerful methods and techniques from several technological domains need to be brought together and integrated. Experience shows that the networks can be unwound and analyzed after the events. Although it provides the necessary evidence for bringing criminals to justice, it is then too late to prevent loss of life and

Memon, N.; Larsen, H.L. (2006) Investigative Data Mining Toolkit: A Software Prototype for Visualizing, Analyzing and Destabilizing Terrorist Networks. In *Visualising Network Information* (pp. 14-1 – 14-24). Meeting Proceedings RTO-MP-IST-063, Paper 14. Neuilly-sur-Seine, France: RTO. Available from: <http://www.rto.nato.int/abstracts.asp>.

Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>01 DEC 2006</b>		2. REPORT TYPE <b>N/A</b>		3. DATES COVERED <b>-</b>	
4. TITLE AND SUBTITLE <b>Investigative Data Mining Toolkit: A Software Prototype for Visualizing, Analyzing and Destabilizing Terrorist Networks</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>Software Intelligence Security Research Center Department of Computer Science and Engineering Aalborg University, Niels Bohrs Vej 8 6700 Esbjerg, Denmark</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release, distribution unlimited</b>					
13. SUPPLEMENTARY NOTES <b>See also ADM002067., The original document contains color images.</b>					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT <b>UU</b>	18. NUMBER OF PAGES <b>68</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

material damage.

Mathematical methods used in our research on Investigative Data Mining [1] [2] are clearly relevant to law enforcement intelligence work and may provide tools to discover terrorist networks in their planning phase and thereby prevent terrorist acts and other large-scale crimes from being carried out. Relevant patterns to investigate include connections between actors (meetings, messages), activities of the involved actors (specialized training, purchasing of equipment) and information gathering (time tables, visiting sites).

Investigative Data Mining (IDM) offers the ability to firstly map a covert cell, and to secondly measure the specific structural and interactional criteria of such a cell. This framework aims to connect the dots between individuals and “map and measure complex, covert, human groups and organisations”. The method focuses on uncovering the patterning of people’s interaction, and correctly interpreting these networks assists “in predicting behaviour and decision-making within the network”.

IDM also endows the analyst the ability to measure the level of covertness and efficiency of the cell as a whole, and also the level of activity, ability to access others, and the level of control over a network each individual possesses. The measurement of these criteria allows specific counter-terrorism applications to be drawn, and assists in the assessment of the most effective methods of disrupting and neutralising a terrorist cell. In short IDM “provides a useful way of structuring knowledge and framing further research. Ideally it can also enhance an analyst’s predictive capability”.

IDM borrowed social network analysis (SNA) and graph theory techniques for connecting the dots; our goal is to propose mathematical methods for destabilizing terrorist networks after linking the dots between them.

In investigative data mining, a number of variations exist in the literature. One is known as link analysis (see for example [5] [6]). Link analysis research uses search and probabilistic approaches to find structural characteristic in the network such as hubs, gatekeepers, pulse-takers [7], or identifying potential relationships for relational data mining. Link analysis alone is insufficient as it looks at one side of the coin and ignores complex nonlinear relationships that may exist between the attributes. Another approach depends purely on visualization, such as NetMap [8]. Unfortunately, these tools that depend on visualization alone - despite being useful to provide some insight - they are insufficient and rely on the user to carry out many tedious and time consuming tasks, many of which could be automated.

In addition to the previous discussion, most of the work on link analysis or network visualization ignores the construction of hidden hierarchy of covert networks. Uncovering a relationship among or within attributes (connecting the dots) is an important step, but in many domains it is more important to understand how this relationship evolved. Hence, understanding network dynamics and evolution is needed to complete the picture. Once we understand the dynamics and evolution of these relationships and construct the hidden hierarchy, we can search for ways to disconnect the dots when and if needed.

The three innovative points of our research are [2]:

1. The use of new measure Position Role Index (PRI) on the pattern of efficiency introduced by Vito Latora and Massimo Marchiori [27]. This measure identifies key players (gatekeepers/ leaders) and followers in the network.
2. The use of another measure known as Dependence Centrality (DC) which discovers who is depending on whom in a network.
3. Estimate possible hierarchical structure of a complex network by applying degree centrality and Eigenvector centrality from social network analysis (SNA) literature and combining it with new measure dependence centrality.

This paper presents some case studies of the terrorist events occurred in the past, using software prototype that we have developed. The structure of the rest of this paper is as follows. In Section 2, a brief introduction about terrorist network analysis is presented. Section 3 presents a concise review to point the reader to several key papers in the literature; whereas Section 4 discusses the models we have used for destabilizing terrorist networks. Section 5 discusses an overview of the iMiner software prototype system; whereas four case studies are then presented in Section 6. Conclusions are then drawn in Section 7.

## 2.0 TERRORIST NETWORK ANALYSIS

The threat from modern terrorism manifests itself worldwide in locally and internationally operating network structures. Before focusing on the development and composition of these networks, we will initially attempt to answer the question: what exactly is a terrorist network?

Persons involved in support, preparation or commission of terrorist attacks almost never operate alone, but as members of - sometimes overlapping - network structures. Within these networks they co-operate with individual members or small groups of members (operational cells). A modern terrorist *network* differs from other terrorist groups and organizations in that it lacks a formal (hierarchical) structure, and has an informal, flexible membership and fluctuating leadership. It is incorrect, however, to conclude that such a network possesses no structure whatsoever. There is always a pattern of connections between individuals who communicate with one another with a view to achieving a common goal. In some cases these communication lines converge in one or more core groups, which thus play a coordinating and controlling role. In other cases there are random communication patterns between all members while the network functions practically without any leadership or central control. It is also possible for several groups to be active within one network.

The flexible and informal character of such a network makes it easy for individual members to establish temporary ad-hoc contacts, in addition to more permanent relations. It also leaves room for personal initiative. The relations within a network are constantly changing in character and duration. In most cases we can distinguish a core group surrounded by a diffuse network of individuals, with central control usually restricted to a minimum. Personal ties between members bind the network together. These relationships are usually based on a shared political-religious ideology, mutual trust, family or friendship ties, shared origin and/or shared experiences in training camps or jihad areas. The notion of a common *enemy* also stimulates bonding among network members.

The above characteristics lead to the following definition:

*A terrorist network is a fluid, dynamic, vaguely delineated structure comprising a number of interrelated persons who are linked both individually and on an aggregate level (cells / groups). They have at least a temporary common interest, i.e. the pursuit of a jihadism-related goal (including terrorism).*

Persons within such a network are referred to as members. A member is a person who contributes actively and consciously to the realization of the aforementioned goal within the bounds of the network.

This definition is in line with the definition of criminal networks used in Criminology, which does not refer to permanent structures, but to temporary, flexible co-operative structures between individuals, based on kinship, friendship, business opportunism, coincidence, necessity, temptation and force, or to the fact that members are colleagues, neighbours or fellow convicts. This co-operation gradually evolves into certain customs and traditions which lead to 'habituation, mutual interdependence and trust, and hierarchical relations [8]. This assessment of fluid and dynamic criminal networks was described in an extensive study into organized crime [9].

### 2.1 Centrality measures for Analyzing Terrorist Networks

Centrality is one of the most important and widely used measures for analyzing social networks. Nearly all empirical studies try to identify the most important actors (also known as vertices / nodes) within the network. Four measures of centrality are commonly used in network analysis: degree, closeness, betweenness, and eigenvector centrality. The first three were described in modern form by Freeman [10] while the last was proposed by Bonacich [11]. Let us begin with degree.

*Degree* centrality measures how active a particular node is. It is defined as the number of direct links a node  $k$  has:

$$C_D(k) = \sum_{i=1}^n a(i, k), \quad (1)$$

where  $n$  is the total number of nodes in a network, and  $a(i, k)$  is a binary variable indicating whether a link exists between  $i^{\text{th}}$  and  $k^{\text{th}}$  nodes. A network member with a high degree could be the leader or “hub” in a network.

*Closeness* centrality [10] is the sum of the length of geodesics (shortest paths between two nodes) between a particular node  $k$  and all other nodes in a network. It actually measures how far away one node is from other nodes and is sometimes called *farness*:

$$C_c(k) = \sum_{i=1}^n l(i, k), \quad (2)$$

where  $l(i, k)$  is the length of the shortest path connecting  $i^{\text{th}}$  and  $k^{\text{th}}$  nodes.

*Betweenness* centrality [10] measures the extent to which a particular node lies between other nodes in a network. The betweenness of  $k^{\text{th}}$  node is defined as the number of geodesics passing through it:

$$C_B(k) = \sum_{i=1}^n \sum_{j=1}^n g_{ij}(k), \quad (3)$$

where  $g_{ij}(k)$  indicates whether the shortest path between two other  $i^{\text{th}}$  and  $j^{\text{th}}$  nodes passes through the  $k^{\text{th}}$  node. A member with high betweenness may act as a gatekeeper or broker in a network for smooth communication or flow of goods (e.g., drugs).

A more sophisticated version of the same idea is the so-called Eigenvector centrality. *Eigenvector centrality*  $x_i$  of a node in a network is defined to be proportional to the sum of the centralities of the node’s neighbours, so that a node can acquire high centrality either by being connected to a lot of others (as with simple degree centrality) or by being connected to others that themselves are highly central. We write

$$x_i = \frac{1}{\lambda} \sum_{j=1}^n A_{i,j} x_j, \quad (4)$$

where  $n$  is total number of nodes and  $\lambda$  is constant. In matrix notation this becomes  $\lambda \mathbf{x} = \mathbf{A} \mathbf{x}$ , so that  $\mathbf{x}$  is an eigenvector of the adjacency matrix [11]. Assuming that we wish the centralities to be non-negative, it can be shown that  $\lambda$  must be the largest Eigen value of the adjacency matrix and  $\mathbf{x}$  the corresponding Eigenvector.

### 2.1.1 Example

Figure 1 shows an example of a terrorist network, which maps the links between terrorists involved in the tragic events of September 11, 2001. This graph was constructed by Valdis Krebs [12] using the public data that were available before, but collected after the event. Even though the information mapped in this network is by no means complete, its analysis may still provide valuable insights into the structure of a terrorist organization. This graph is reconstructed in this paper, using metadata (additional information) of every terrorist involved in the attacks.

According to Krebs’s analysis [12], this network had 62 members in total, of which 19 were kidnapers, and 43 assistants: organizers, couriers, financiers, scouts, representatives, coordinators, counterfeiters, etc. Allen [13] found that successfully functioning large networks typically comprise 25-80 members, with optimal size between 45 and 50. A close match exists between the results of Allen’s analysis of collaborating networked groups and this particular example of terrorist group.

Inspection of this network by standard measures of network structure reveals firstly its low connectedness. A member of this network holds only 4.9 connections with others members on average (also known as degree centrality), which means that average members were rather isolated from the rest of the network. The density (which is defined as the number of actual links divided by the number of possible links) of



this network is only 0.08, meaning that only 8% of all possible connections in the network really exist (see Figure 2).

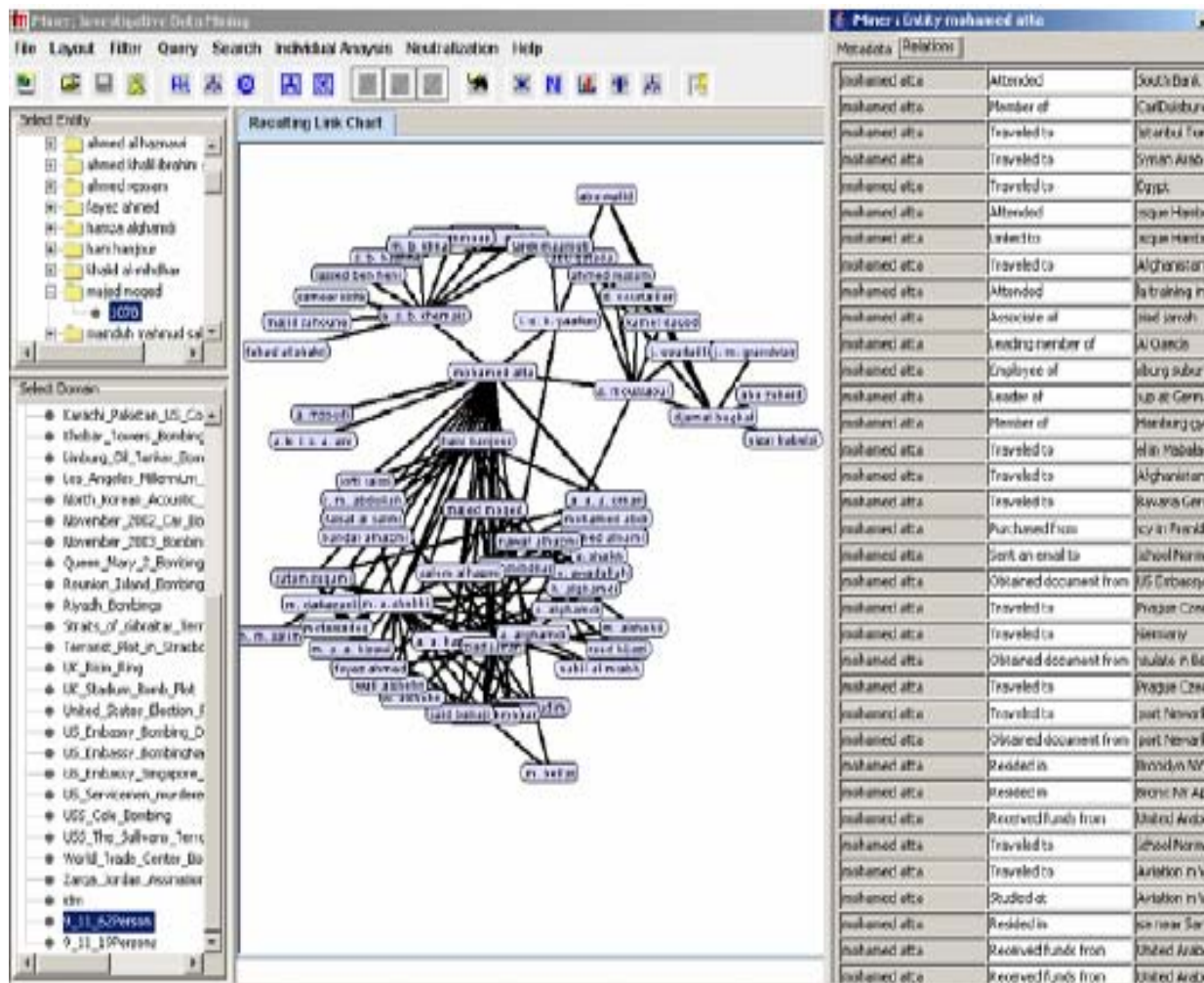


Figure 1. The dataset of 9-11 hijackers and their affiliates. The dataset originally constructed by Valdis Krebs, but re-constructed in iMiner, using additional information of every entity.

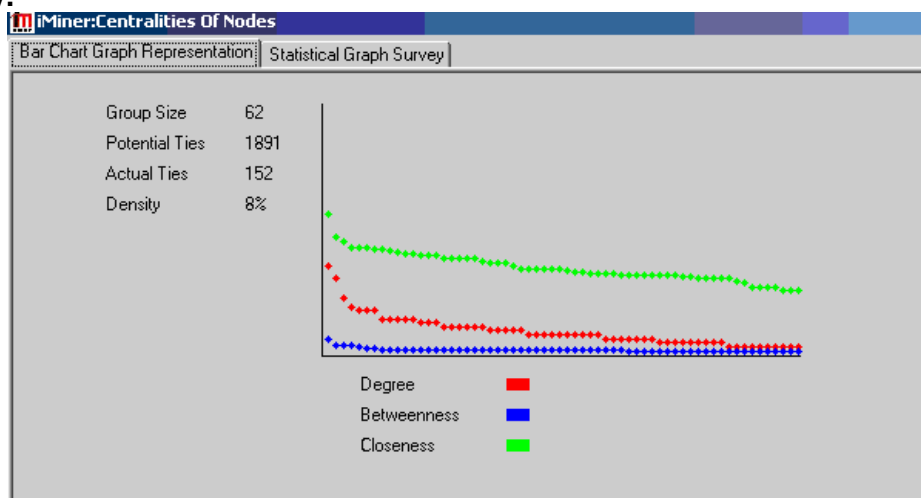
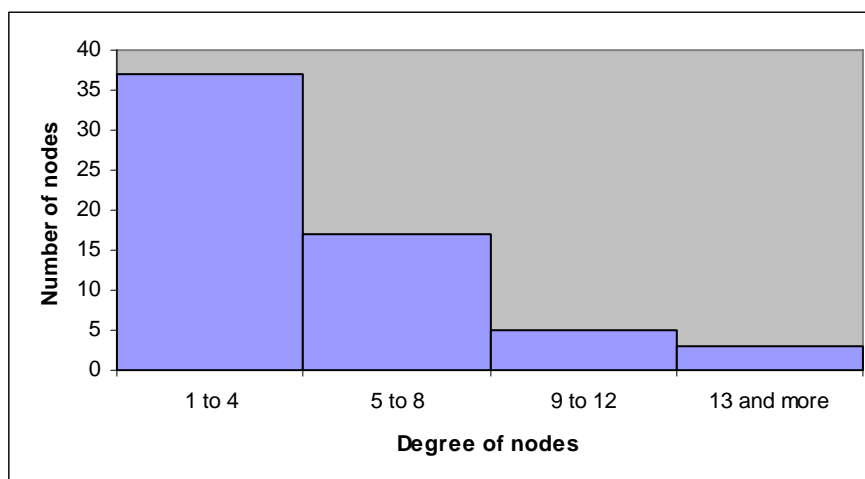


Figure 2. 9-11 Terrorists Network's Neighbourhood

In spite of low connectedness, however, the nodes of this network are relatively close. The average closeness of nodes is 0.35. Betweenness as stated above is another important measure in SNA and it indicates a node's importance for communication among other nodes. The average betweenness of this network is 0.032, indicating relatively high average redundancy. However the betweenness of 40 nodes is in fact less than 1% and only 6 nodes have betweenness higher than 10%. These 6 nodes are critical for information flow, especially one with betweenness of almost 60%, meaning that almost 60% of communication paths among other nodes pass through this central node. The node represents *Mohamed Atta*, the leading organizer of the attack whose central position in the network is confirmed by other centrality indicators as well.

Distribution of degrees of nodes is particularly interesting. Degrees of nodes are exponentially distributed: the degree of most of the nodes is small, while few nodes have high degree (see Figure 2 and Figure 3). This property characterises the so called scale free networks [14] [15]. Scale free networks form spontaneously, without needing a particular plan or interventions of central authority. Nodes that are members of the network for a longer time, that are better connected with other nodes, and that are more significant for network's functioning, are also more visible to new members, so that the new members spontaneously connect more readily to such nodes than other, relatively marginal ones.



**Figure 3. Distribution of degrees of nodes in the network (see Figure 1) of kidnappers and their supporters.**

On the pattern of scale free networks, the Al Qaeda's Training Manual [16] states: "The cell or cluster methods should be organized in a way that a group is composed of many cells whose members do not know each other, so that if a cell member is caught, other cells would not be affected, and work would proceed normally".

## **2.2 Link Analysis**

Link analysis is an analytic technique for making relationship explicit. The link analysis is the process of building up networks of interconnected objects through relationship in order to expose patterns and trends. Link analysis uses item-to-item associations to generate networks of interactions and connections from defined datasets. Link analysis methods add dimensions to an analysis that other forms of visualization do not support. Link analysis can be used to construct inferential structures of organizations or interactions which can be tested later. It is very well suited for hypothesis construction and can be applied to a variety of problems in Armed Forces Intelligence (for example, orders of battle), in Political Intelligence and in

Sociological Intelligence research and analysis [17].

Despite the seeming novelty of link analysis, the federal government in USA has used link analysis, for nearly fifty years. Karl Van Meter describes the two main types of the link analysis: the village survey method and traffic analysis [29]. The village survey method was created and used by Ralph McGehee of the CIA in Thailand in the 1960s to understand family and community relationships. He conducted a series of open-ended interviews and in a short time was able to map out clandestine structure of local and regional Communist organizations and associated sympathetic groups. Traffic analysis (also known as communication link analysis) began during World War II and its importance continues to this day. This technique consists of the study the external characteristics of communication in order to get information about the organization of the communication system. It is not concerned with the content of phone calls, but is interested in who calls whom and the network members, messengers, and brokers. Traffic analysis was also used by the British MI5 internal security service to combat the IRA in the 1980s and 1990s and continued to be used across the world by law-enforcement agencies including the US Defense Intelligence Agency (DIA) Office of National Drug Control Policy [29].

The *Analyst Notebook* [32] is the primary software used for link analysis. Currently on its sixth version, this software is recognized as one of the world's leading analytical tools and is employed in more than 1,500 organizations. SNA improves upon link analysis by moving from single variable analysis to multivariate analysis, allowing the individual to control many factors at once. The change from single variable to multivariate analysis is quite significant when researching terrorism: a number of factors affect terrorism, not one single factor. For example, the prediction for one to participate in a terrorist activity might not be strongly affected by the single variable of being related to a terrorist member. However, the combination of multiple variables, such as poverty and type of government, combined with the link to a terrorist member, may cause a person to participate in a terrorist activity. Multivariate analysis allows us to take into account these multiple variables and their effects when controlling another variable.

From the outside, it is difficult to understand how link analysis is being used in the federal governments. Confidentiality prevents government analysts from discussing their work with researchers and private companies without security clearances. Despite this lack of information, it is clear that government is interested in using network techniques in dealing with counterterrorism strategies. Many government agencies, such as Defense Advanced Research Projects Agency (DARPA), U.S. Army Research Labs, the U.S. Office of Naval Research (ONR), the National Security Agency (NSA), the National Science Foundation (NSF), and the Department of Homeland Security (DHS) have funded research related to link analysis.

The link analysis systems are being used in the investigative world, but from unclassified work, it is known that they are only used for identifying central players and some interesting patterns from the available datasets. Apparently little work is carried out for destabilizing of terrorist networks [1] [2] [4] [18] [26]. The motivation behind this study is to connect the dots in order to assist law enforcement agencies to disconnect / destabilize the most of the network by capturing/eradicating some key players.

### **3.0 EXISTING APPROACHES**

Existing terrorist network research is still at its incipient stage. Although previous research, including a few empirical ones, have motivated the call for new approaches to terrorist network analysis [18] [19] [20], studies have remained mostly small-scale and used manual analysis of a specific terrorist organization. For instance, Krebs [12] manually collected data from public news releases after the 9/11 attacks and studied the network surrounding the 19 hijackers. Sageman [21] analyzed the Global Salafi Jihad network consisting of 171 members using a manual approach and provided an anecdotal explanation of the formation and evolution of this network. None of these studies used advanced data mining technologies that have been applied widely in other domains such as finance, marketing, and business to discover previously unknown patterns from terrorist networks.

The papers [18] [22] provide examples of social network analysis in counterterrorism applications and



indicate both usefulness and some limitations of social network analysis as a basis for quantitative methods for situation awareness and decision-making in law enforcement applications. Raab & Milward [22] discuss the organizational structure of certain drug trafficking, terrorism, and arms trafficking networks, showing how some of them have adapted to increased pressure from the States and international organizations by decentralizing into smaller units linked only by function, information, and immediate need. They also describe ways and structures of cooperation between different kinds of criminal networks, for example in financing terrorists by illegal diamond and drug trafficking.

Recently, computer scientists have become interested in network analysis. This has led to an increased emphasis on studying the statistical properties of large networks, such as the Internet, criminal networks, and, even, infrastructure networks [23]. This influx of people to the field has also led to several new approximate algorithms to compute important properties [1] [2] [24] [25]. Most of the above mentioned literature has used graph theory and SNA techniques for terrorist network analysis.

Jonathan D. Farley presented a new mathematical approach [26] to destabilize terrorist networks using order theory. This paper pointed that modeling terrorist networks as graphs does not give enough information to deal with the threat, "modeling terrorist cells as graphs ignores an important aspect of their structure, namely their hierarchy, and the fact that they are composed of leaders and followers". Jonathan D. Farley proposed an alternative approach that better reflects an organization's hierarchy. In this case, the relationship of one individual to another in a cell becomes important. Leaders are represented by the topmost nodes in a diagram of the ordered set representing a cell and foot soldiers are nodes at the bottom. Disrupting the organization would be equivalent to disrupting the chain of command, which allows orders to pass from leaders to foot soldiers [26].

What is needed is a set of integrated methods, technologies, models, and tools to automatically mine data and discover valuable knowledge from terrorist networks based on large volumes of data of high complexity.

The system to be proposed, we represent a terrorist network by an undirected graph; then we convert it into a directed graph with the help of centrality measures [1] [2]. Then we propose an approach for destabilizing the terrorist network to the converted directed graph into a hierarchical chart using dependence centrality [1]. From the hierarchical chart, investigators and law enforcement agencies can easily distinguish the leaders and peripheries in the network in order to destabilize the network. Our newly introduced dependence centrality measure may also be very useful in destabilizing terrorist networks, because it shows the nodes which are totally depending on particular nodes. If the nodes are completely depending on the other nodes, they will be isolated (cut-off from the network completely) by capturing the node on which those nodes are depending.

## **4.0 DESTABILIZING TERRORIST NETWORKS**

In this section we present a theory behind analyzing and destabilizing of terrorist networks. We have implemented all the models discussed in this section in iMiner.

### **4.1 The Efficiency of a network**

The network efficiency  $E(G)$  is a measure to quantify how efficiently the nodes of the network exchange information [27]. To define efficiency of a network  $G$ , first we calculate the shortest path lengths  $d_{ij}$  between  $i^{\text{th}}$  and  $j^{\text{th}}$  nodes. Let us now suppose that every node sends information along the network, through its links. The efficiency in the communication between  $i^{\text{th}}$  node and  $j^{\text{th}}$  node is inversely proportional to the shortest distance: when there is no path in the graph between  $i^{\text{th}}$  and  $j^{\text{th}}$  nodes, we get  $d_{ij} = +\infty$  and efficiency becomes zero. Let  $N$  is known as the size of the network or the numbers of nodes in the graph, the average efficiency of the graph (network) of  $G$  can be defined as:

$$E(G) = \frac{1}{N(N-1)} \sum_{i \neq j \in G} \frac{1}{d_{ij}} \quad (5)$$

The above formula gives a value of E in the interval of [0, 1].

## 4.2 Position Role Index (PRI)

The PRI is our recently introduced measure [2] which highlights a clear distinction between followers and gatekeepers (It is a fact that leaders may act as gatekeepers). It depends on the basic definition of efficiency as discussed in equation (5). It is also a fact that the efficiency of a network in presence of followers is low in comparison to their absence in the network. This is because they are usually less connected nodes and their presence increases the number of low connected nodes in a network, thus decreasing its efficiency.

$$PRI = E(G) - E(G - v_i), i = 1, \dots, N, \quad (6)$$

where  $G - v_i$  indicates the network obtained by deactivating node  $v_i$  in the graph G. If we plot the values on the graph, the nodes which are plotted below x-axis are followers, whereas the nodes higher than remaining nodes with higher values on positive y axis are the gatekeepers.

## 4.4 Dependence Centrality (DC)

The dependence centrality of a node is defined as how much that node is dependent on any other node in the network. It is defined as:

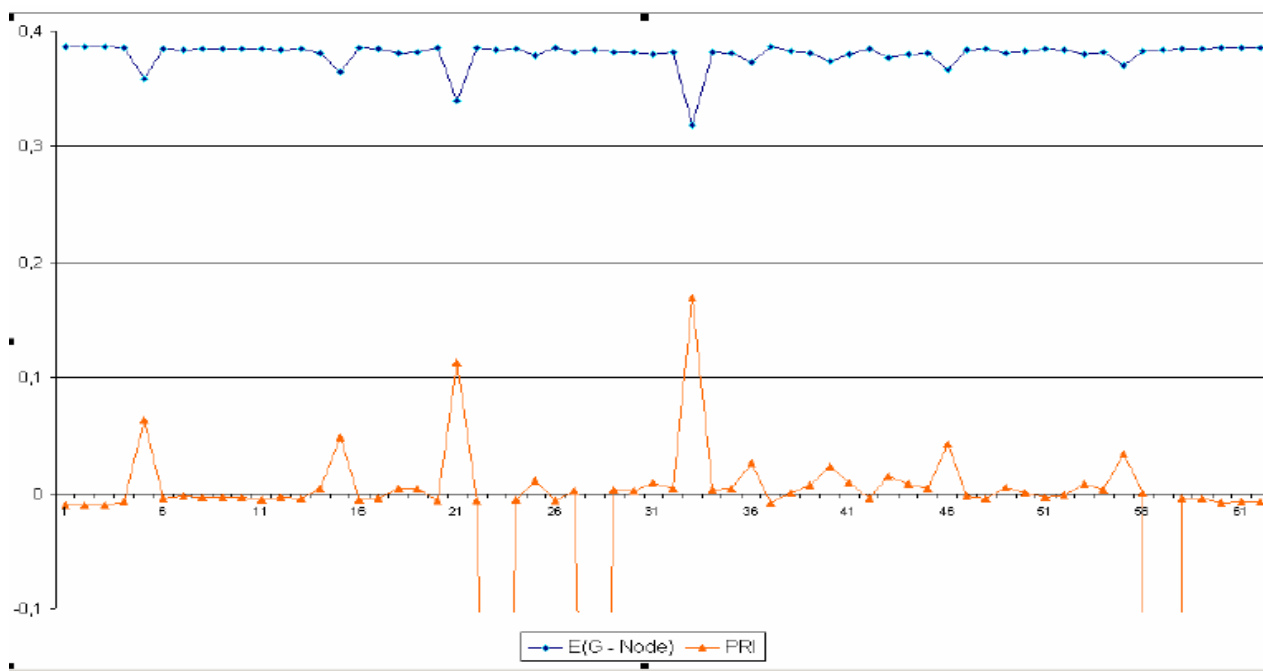
$$C_\delta(i) = \sum_{i \neq k, k \in G} \frac{m(i, j)}{N_p} + \Omega, \quad (7)$$

where  $i^{\text{th}}$  node is the root node which depends on  $j^{\text{th}}$  node,  $N_p$  is the number of geodesic (shortest path) from  $i^{\text{th}}$  node to node  $k^{\text{th}}$  through  $j^{\text{th}}$  node. Also  $m(i, j)$  is inverse of geodesic distance  $1/d(i, j)$  from  $i^{\text{th}}$  node to  $j^{\text{th}}$  node. The value of  $\Omega$  is taken 1 if graph is connected and 0 in case it is disconnected. In this paper we take  $\Omega$  as 1, because we consider that graph is connected.

The first part of the formula tells us that:

How many times  $i^{\text{th}}$  node uses  $j^{\text{th}}$  node to communicate with other  $k^{\text{th}}$  node of the network? In simple words  $k$  is the node of the network, to which node  $i$  is connected through node  $j$ . (The connection represents the shortest path of node  $i$  to node  $k$ , and node  $j$  is in between).

We applied the above mentioned measures (described in subsection 4.1-4.2) in the network of terrorists involved in tragic events of September 11, 2001 (as shown in Figure 1). The results are depicted in Figure 4. The results show that node 33 (Mohamed Atta) as key player in the plot. The  $\Delta E$  and PRI of this node is higher than all nodes which prove that this node played an important role in the plot and worked as gatekeeper and removing this node the efficiency of the graph is decreased from 0.395 to 0.32. This clearly identifies the importance of this node in the network. All the models presented in previous sections are implemented in iMiner.



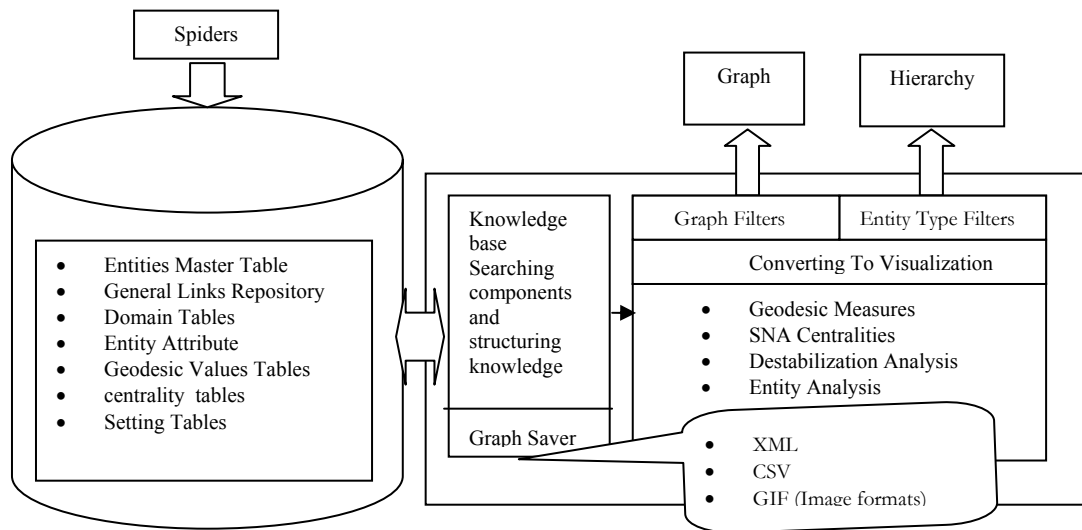
**Figure 4.** The efficiency of the original network  $E(G) = 0.395$ . The removed node is shown on x-axis. The efficiency of the graph once the node is removed is shown as  $E(G - v_i)$ ; while importance of node. The newly introduced measure position role index is shown as PRI.

## 5.0 THE SOFTWARE PROTOTYPE

The iMiner is an experimental system, which provides the answers of the following questions.

- Who is likely to be an important person (node) in the network?
- Why s/he is important?
- Which terrorist is highly/ less connected?
- What are various position roles in the network?
- Is there any command structure in the network?
- Is it possible to construct hidden hierarchy of the non-hierarchical networks?
- Which nodes represent key players?
- What is the efficiency of the network?
- How much the efficiency of the network is affected by eradicating one or more terrorists?
- How can the law enforcement agencies use (often incomplete and faulty) network data to disrupt and destabilize terrorist networks?

The system architecture of investigative data mining toolkit (iMiner prototype) is depicted in Figure 5. The first stage of network analysis development is intended to automatically identify the strongest association paths, or geodesics, between two or more network members. In practice, such task often entails intelligence officials to manually explore links and try to find association paths that might be useful for generating investigative leads.



**Figure 5. The System Architecture of iMiner**

The iMiner allows the analyst to determine hierarchy of covert networks, which may help law enforcement and intelligence agencies in understanding the structure of the covert networks. In addition to this, intelligence agencies can easily set the priorities for eradicating some important nodes in the network by visualizing how much the efficiency of the network is minimized by the capture of a node.

The iMiner knowledgebase has following type of tables

- Entities Master Table
- General Links Repository
- Domain Tables
- Entity Attribute
- Geodesic Values Tables
- centrality tables
- Setting Tables

The above tables are briefly described below:

#### **Entities Master Table**

This table is a table where entity is assigned with unique id and very primary information about entity like entity's name and type being saved.

#### **General Links Repository**

This table has a record of all relationships in form of Association Data Model. These all relationships are used while searching the database and graphing a particular plot.

#### **Domain Tables**

These tables are subsets of general links repository. The relationships which are connected to a particular domain are stored separately in that domain table so that the complex operations which performed by iMiner during analysis takes place efficiently.

### **Entity Attribute Tables**

These tables are used to store attributes of entities.

### **Geodesic Values Tables**

These tables are used to save the geodesic values which are calculated during analysis of a particular domain and are useful for other analysis activities. The iMiner has components building to show those values in anytime to help analysts.

### **Centrality Tables**

These tables are used to save the values of centralities (SNA centralities, dependence centrality, position role index, etc.) of each node in a graph of a particular domain under investigation.

### **Setting Tables**

There are some tables which are used to save settings of iMiner, like search depth, entity types along with their visualization scheme, etc.

In the aftermath of the September 11<sup>th</sup> attacks, it was noted that coherent information sources were not available to the researchers [3]. Information was either available in fragmentary form, not allowing comparison studies across incidents, groups or tactics, or made available in written articles – which are not readily suitable for quantitative analysis of terrorist networks. Data collected by law enforcement agencies, while potentially better organized, are largely not available to the research community due to restrictions in distribution of sensitive information.

To counter the information scarcity, we have developed knowledgebase about the terrorist attacks occurred in the past and the information about terrorist organizations involved in those events. This information is mostly collected from open source media (but authenticated websites), such as <http://www.trackingthethreat.com/>.

The focus of the knowledgebase we have developed is the agglomeration of publicly available data and integration of the knowledgebase with investigative data mining software prototype. The main objective is to investigate and analyze terrorist networks to find hidden relations and groups, prune datasets to locate regions of interest, find key players, characterize the structure, trace point of vulnerability, and detect efficiency of the network and to discover the hidden hierarchy of the non-hierarchical networks. The Website <http://www.trackingthethreat.com/> represents the original open-source database on Al Qaeda. It contains data in the form of:

**Entities:** Discrete data elements that comprise people, places, organizations, events, etc.

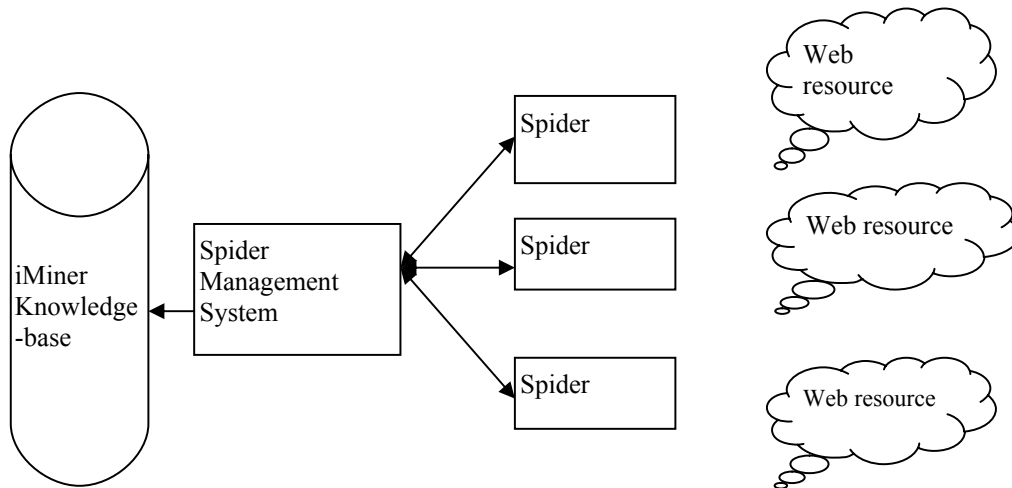
**Relationships:** Information about the personal, organizational, transactional, and historical connections between entities.

**Metadata:** Additional information about entities and relationships that help form a more complete picture.

**Notes and Documents:** Unstructured text that provides background information on entities, relationships, and metadata.

The iMiner system applies Spider Management System as well as spiders to import data that is available at World Wide Web. It can be developed to get information from online repositories and save it in its knowledgebase for analysis as shown in Figure 6.



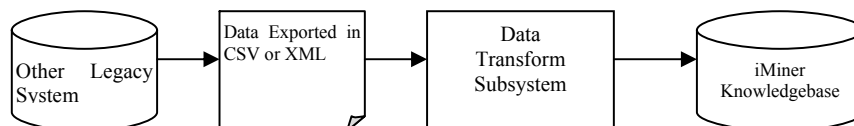


**Figure 6. Spidering the information from Internet**

The iMiner spidering system has spider implementation for a web resource as mentioned above (<http://www.trackingthethreat.com/>). Similarly this model can be evolved to make most of information present in shape of static and dynamic websites. The data flow model for the prototype is shown in Figure 8.

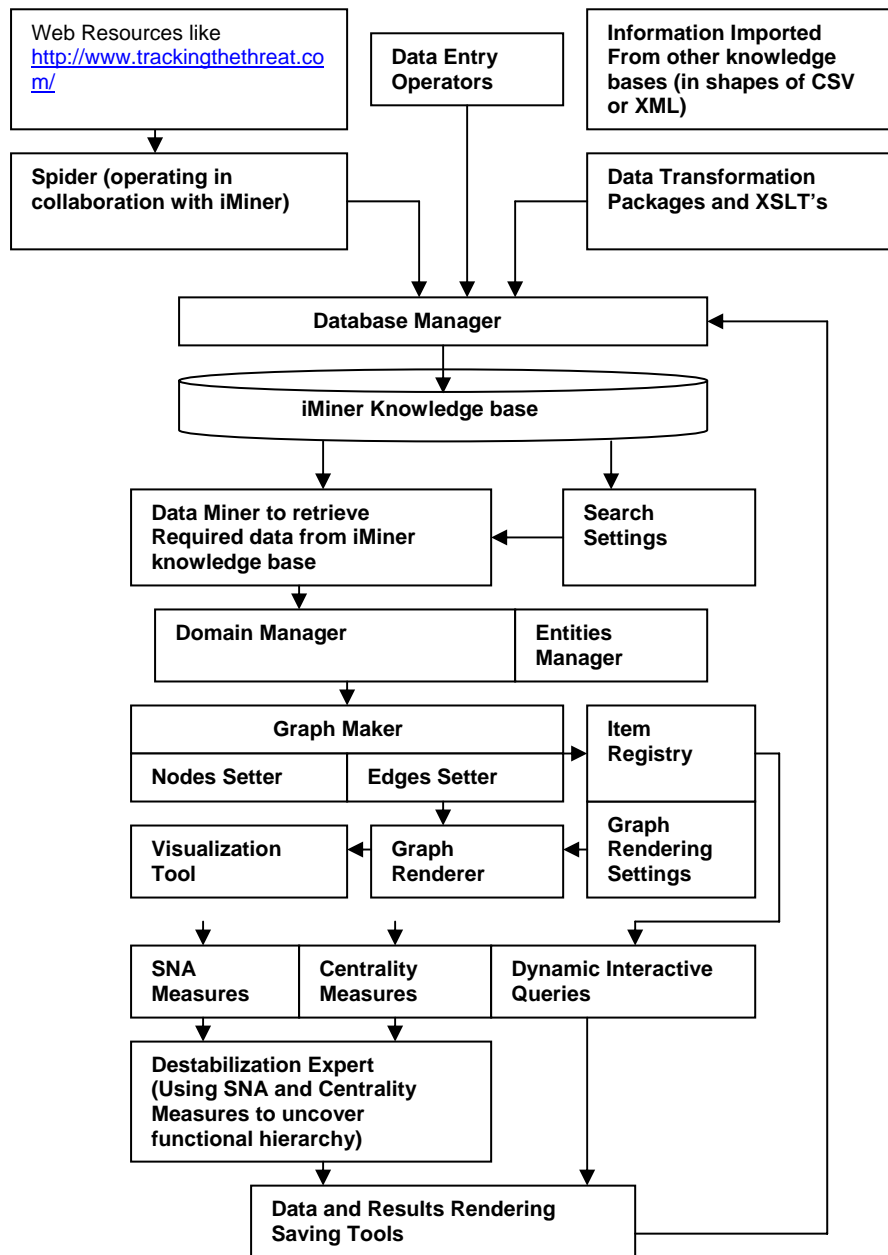
The iMiner supports domain-based study of network system. The domain can be similar to a particular case study (for example Bali Bombing or 9/11 case studies) with resizable boundaries. An investigating officer may expand the domain by including additional entities or reduce it by excluding less important or partially involved entities. The main advantage of this type of approach is to isolate a piece of data in which an investigating officer is interested, to concentrate analysis activities only on those entities.

To reuse the type of data, the iMiner has a data transform subsystem; the purpose of the subsystem is to import data from different legacy systems and transform them into the format required by iMiner knowledge base to carry out further investigation as shown in Figure 7.



**Figure 7. The data transform sub-system**

Most of the commercial database solutions provide facilities to export data in at least CSV or XML format. The information exported from any legacy system or data source in CSV or XML format can be imported to extend iMiner knowledgebase to proceed with further investigation.



**Figure 8. The data flow model of iMiner**

## 6.0 CASE STUDIES

In this section, we present four case studies of the terrorist attacks occurred / planned in the past to test the algorithms and models recently published [1] [2].

### 6.1 Bali Night Club Bombing Terrorist Attack

On October 12, 2002 a car bomb exploded outside the Sari night club in Bali, a popular tourist island in Indonesia [28]. The attack was the worst terrorist incident in the history of Indonesia, with 202 civilians dead and more than 100 wounded. According to the Australian Broadcasting Corporation, the mastermind,

behind the attacks was Al-Qaeda's chief representative and senior planner in Southeast Asia, as well as being operational chief of Jemaah Islamiya, Ryuduan bin Isomuddin, also known as Hambali. Hambali was detained by the U.S. government in August of 2003. Hambali was also believed to have been involved in the 2003 Marriot Hotel bombings in Jakarta, facilitated the January 2000 meeting in Malaysia including two September 11th hijackers, and an associate of 9/11 mastermind Khalid Sheikh Mohammed.

The Bali Night Club Bombing 2002, when searched in iMiner knowledgebase, results the graph as shown in Figure 9. Using algorithms for destabilizing terrorist networks (for details see [2]), we succeeded in constructing the hidden hierarchy of Bali bombing terrorist attack, shown in Figure 10, using iMiner. This hierarchy has some unconnected nodes, where as you can find a hint of patterns some time. The H. B. A. Haq and its descendants form a group (This cluster was acted as executive cluster), while the cluster Khalid Sheikh Mohammed (and his affiliates) was well known as a strategic cluster, whereas Ryuduan bin Isomuddin (known as Hambali) and his associates cluster known as tactical / logistic cluster. The accuracy of the software can be determined by the fact that all of H. B. A. Haq, Khalid Sheikh Mohammed and R. Isomuddin were key players in the reality. H. B. A. Haq was termed as potential leader while Khalid Sheikh Mohammed was the key conspirator.

In this case study we consider the connections network of terrorist involved in Bali Night Club Bombing and their direct or indirect relationships with other entities. Of course, mapping networks after an event is relatively easy, while the real problem in this case is to map the covert networks to prevent terrorist activity, a task that can be more difficult. The network reported in Figure 9 is constructed by iMiner, using publicly released information taken from major newspapers/ websites. The network size (say N) is 125 nodes, actual ties (arcs/ links) are 195, potential ties, which can be calculated by using the formula:

$N(N-1) / 2$ , are 7750. The density of the network (the number of actual ties divided by the number of potential ties) is approximately 2.5%. To individuate the critical nodes, *i.e.*, the terrorists who played key role in the network, we deactivate one by one (from the network, remember the size of the network is 125), then we calculate the efficiency of the new network and the drop of efficiency caused, discussed in section 4. In Table 1 we show 6 most important nodes ranked according to the measure position role index defined in section 4.2, the degree of each node (represented as k) is also reported.

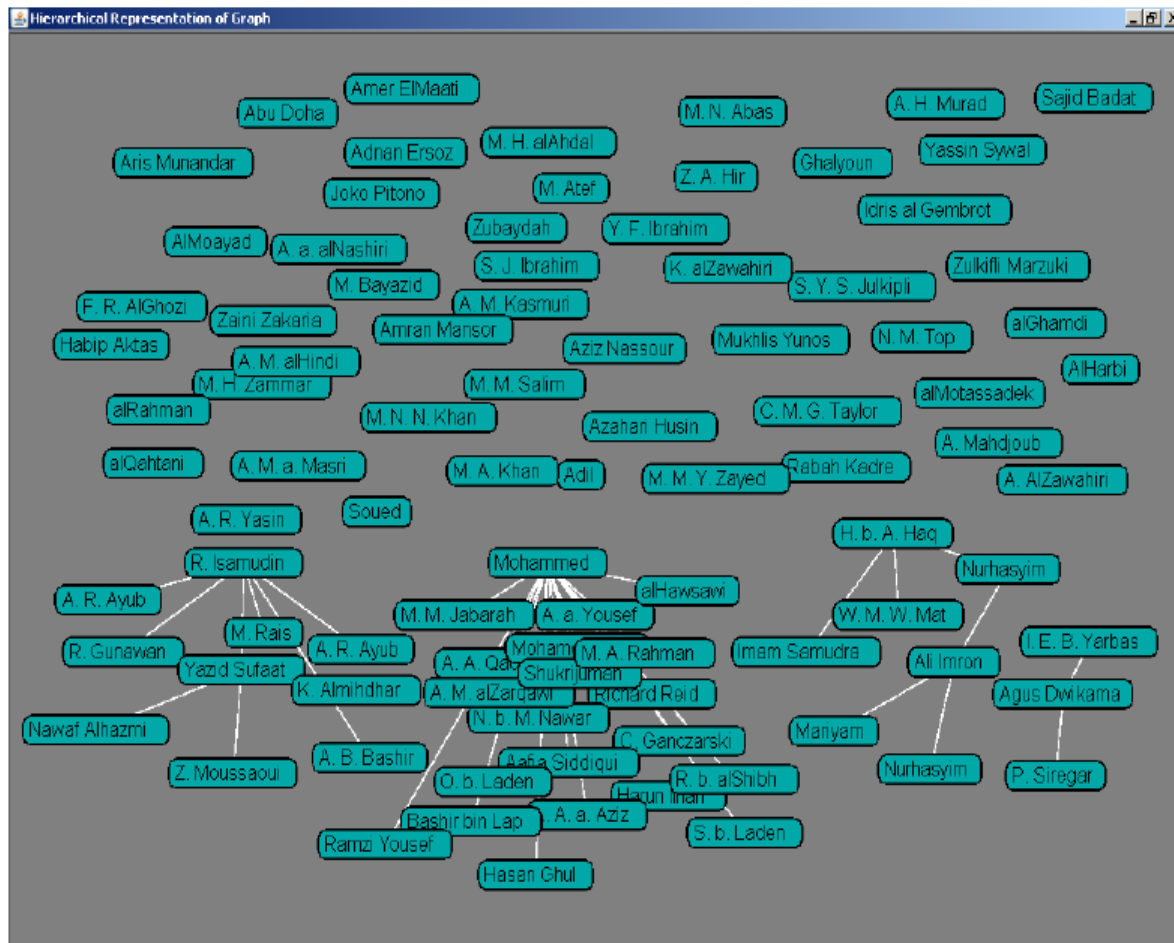
**Table 1. Effect of the deactivation of a node in the terrorist network of Bali Night Club Bombing 2002 attacks**

<b>Removed Node</b>	<b>E(G-v<sub>i</sub>)</b>	<b>k</b>	<b>PRI</b>
K. S. Mohammed	0.3063	27	0.1854
Riduan Isamudin	0.3329	23	0.1148
Yazid Sufaat	0.3342	12	0.1112
Wan Min	0.3700	11	0.0161
Huda Bin A. Haq	0.3745	12	0.0042
Osama Bin Laden	0.3745	3	0.0027



It is crystal clear from the Table 1, that *K.S. Mohammed* was the most connected node in the network ( $k=27$ ), and our recently introduced measure i.e. PRI also proves him as an important entity (*key player*), removing this node the efficiency of the network is decreased to 0.306. Moreover the measure PRI of *Huda Bin A. Haq* is lower than other leaders of two clusters, shows that he was the leader execution cluster. It is also interesting to note that in this network *Osama Bin Laden* is less connected node ( $k=3$ ) and position role index measure shows his leading capabilities. Moreover, 64% of the entities in this network (80 out of 125) having degree ( $k$ ) is equal to 1 (it means about  $2/3^{\text{rd}}$  of the network members degree is 1), entity Al Qaeda is the most connected node having  $k = 52$  (it means that 52 nodes have directly affiliated with the terrorist organization Al Qaeda) and 29 nodes are directly connected with Bali

bombing terrorist plot.



**Figure 10. Hidden hierarchy constructed by iMiner for Bali Bombing Attack**

## 6.2 Dirty Bomb Plot

Adnan Gulshair el Shukrijumah is a suspected of plotting to carry out a terrorist attack against U.S. interests abroad or domestically [30]. When Khalid Sheikh Mohammed, former operational commander of Al Qaeda, was captured and interrogated, he fingered Adnan el Shukrijumah as the man who would later be in charge of new attack. Shukrijumah is believed to be working on Osama bin Laden's plan to trigger a radiological disaster inside the United States [33] – the so-called “dirty-bomb” scenario where a small charge would trigger dispersion of radiation over a large area, wreaking havoc on those caught in the blast and making the blast area uninhabitable. High-grade uranium is not necessary for this project; ordinary, low-grade nuclear waste will be deadly enough.

El Shukrijumah has eluded capture. According to an FBI informant, El Shukrijumah was spotted 2002 in Hamilton, Ontario, posing as student at McMaster University, which has 5-megawatt research reactor. The U.S. officials believe El Shukrijumah was in Hamilton to obtain radioactive material. We have collected dataset for dirty bomb terrorist network plot which is depicted in Figure 11.



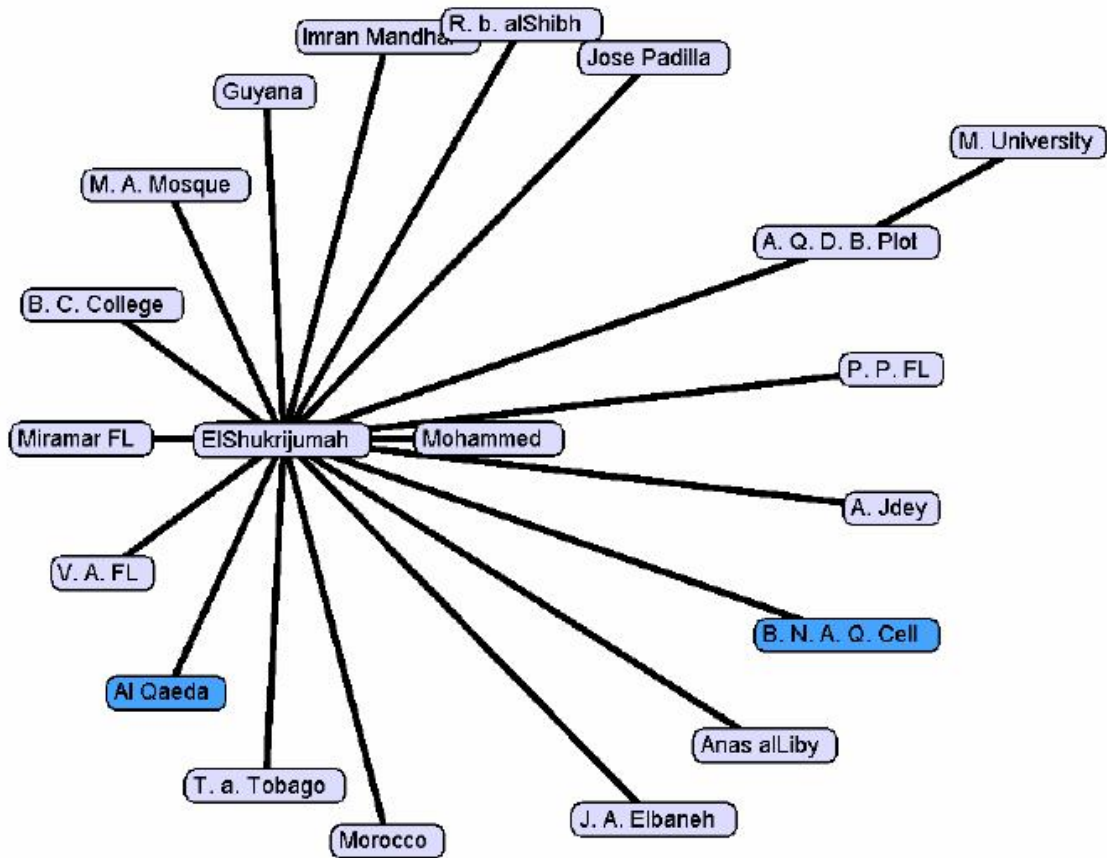


Figure 11. Terrorists' involvement for Dirty Bomb plot

Using the algorithms for construction of hidden hierarchy [2], the hierarchy for the dirty-bomb plot is shown in Figure 12.

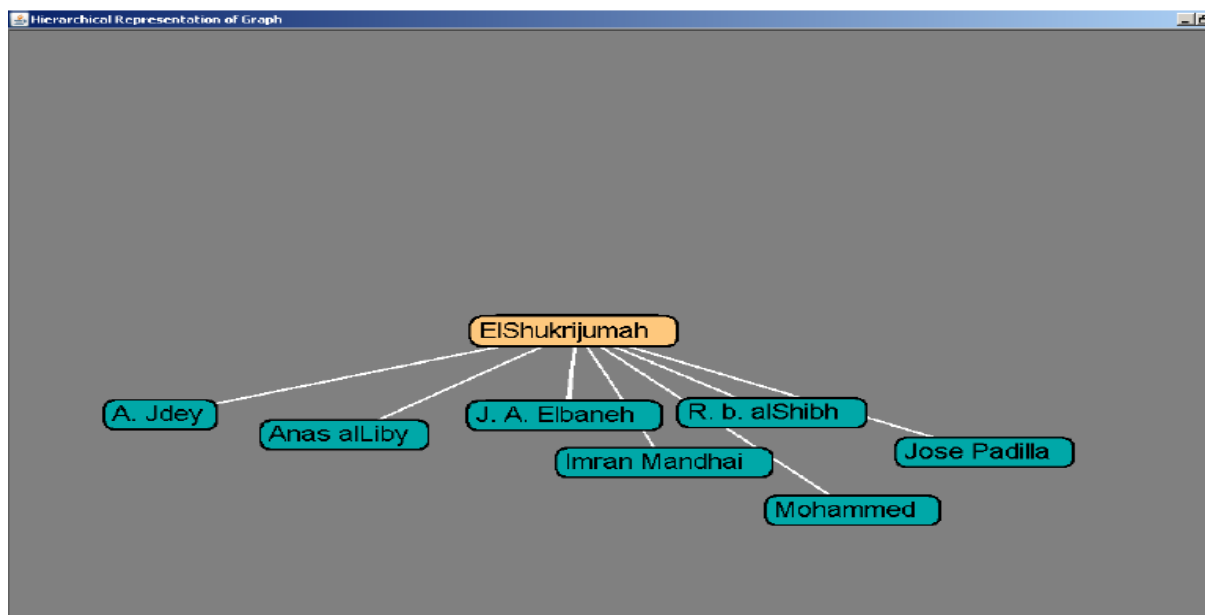


Figure 12. Hidden hierarchy constructed by iMiner for Dirty-Bomb Plot

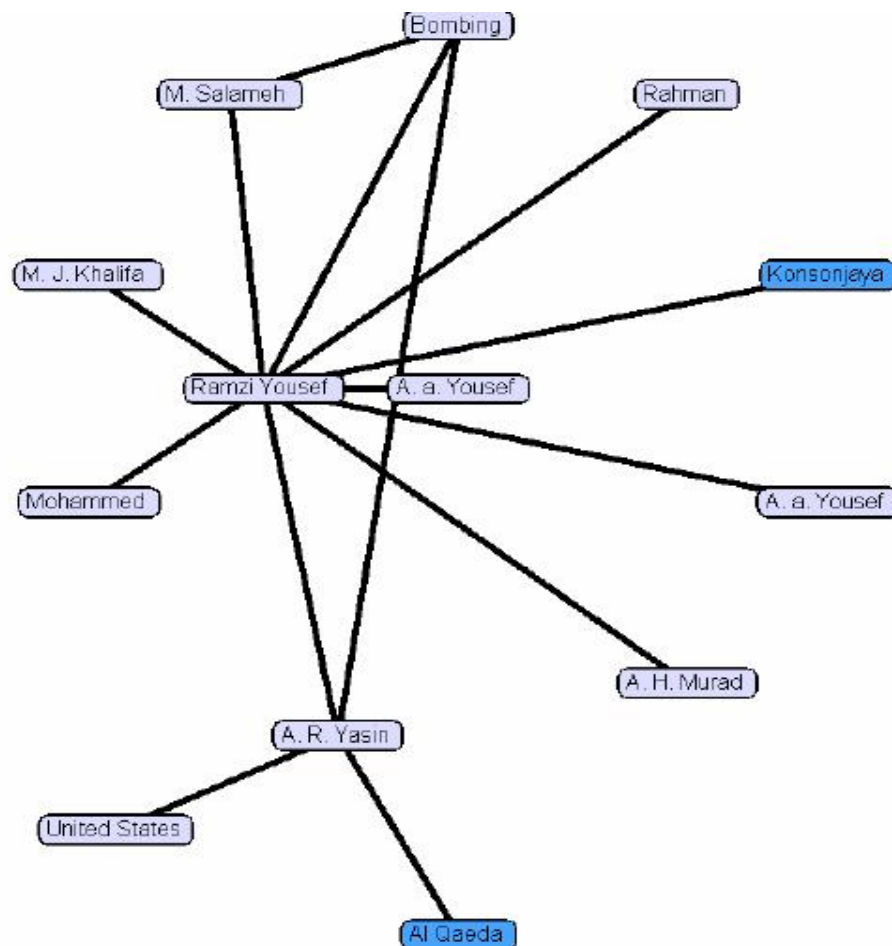
It is believed by authorities that Shukrijumah may have been trained at an Al Qaeda terrorist camp[34]. Shukrijumah has extensive flight training that he received at a flight school in Florida and is a pilot, though he is not registered with the Federal Aviation Administration (FAA). Obviously, this is of concern to law enforcement since virtually all the Al Qaeda hijackers involved in the September 11 terrorist attacks, received training to be pilots at U.S. private flight schools.

Also, it is believed by the FBI that Shukrijumah has been trained by Al Qaeda to operate as a terrorist organizer and operational / field commander and lead or coordinate a terrorist assault, much the same way Mohammed Atta was designated and trained as an organizer and operational/field leader by Al Qaeda to lead the 9/11. El Shukrijumah may play a large and leading role in the next set of terrorist attacks to come upon the U.S. Shukrijumah was last seen in the Miami or southern Florida area in the early part of 2003. He has not been seen since and no one knows of his whereabouts [34].

### 6.3 WTC 1993 Bombing Plot

The WTC bombing attack in the garage was occurred on February 26, 1993 of the New York World Trade Center. A car bomb was planted by terrorist group in the underground parking garage below tower One. It killed six, injured over 1,000 and presaged the 9/11 attacks on the same buildings.

We have collected data of the terrorists involved in the attack and the network is shown in Figure 13.

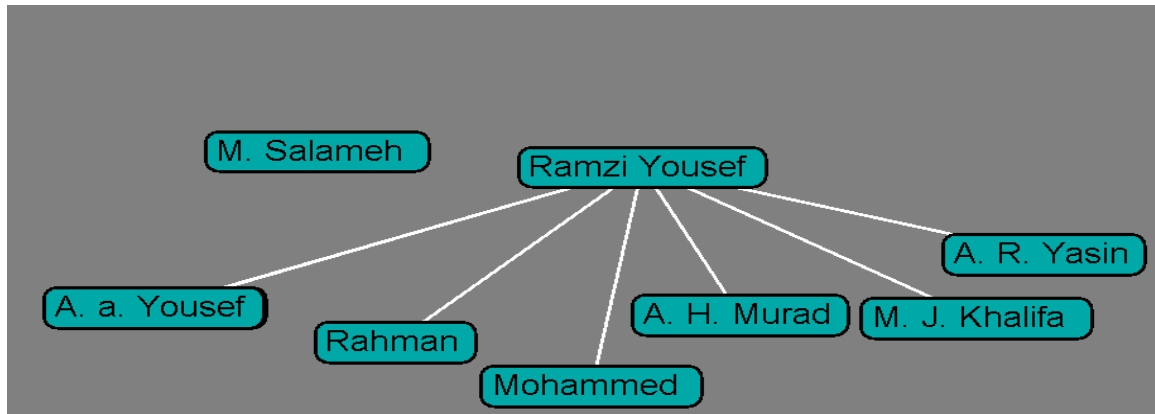


**Figure 13. Terrorists Network involved in WTC 1993**

The group size of the network is 13; potential ties are 78; whereas actual ties are 14. The density of the

network is 18% and efficiency of the network is 0.555. The most important node in the network is Ramzi Yousef ( $k=10$  and Eigenvector centrality = 16). Removing this node the efficiency of the graph is decreased to 0.0845.

Using algorithms for the construction of hidden hierarchy of covert/ terrorist networks [2], the hierarchy of the network is shown in Figure 14.



**Figure 14. Hidden hierarchy constructed by iMiner for the terrorist network shown in Figure 13**

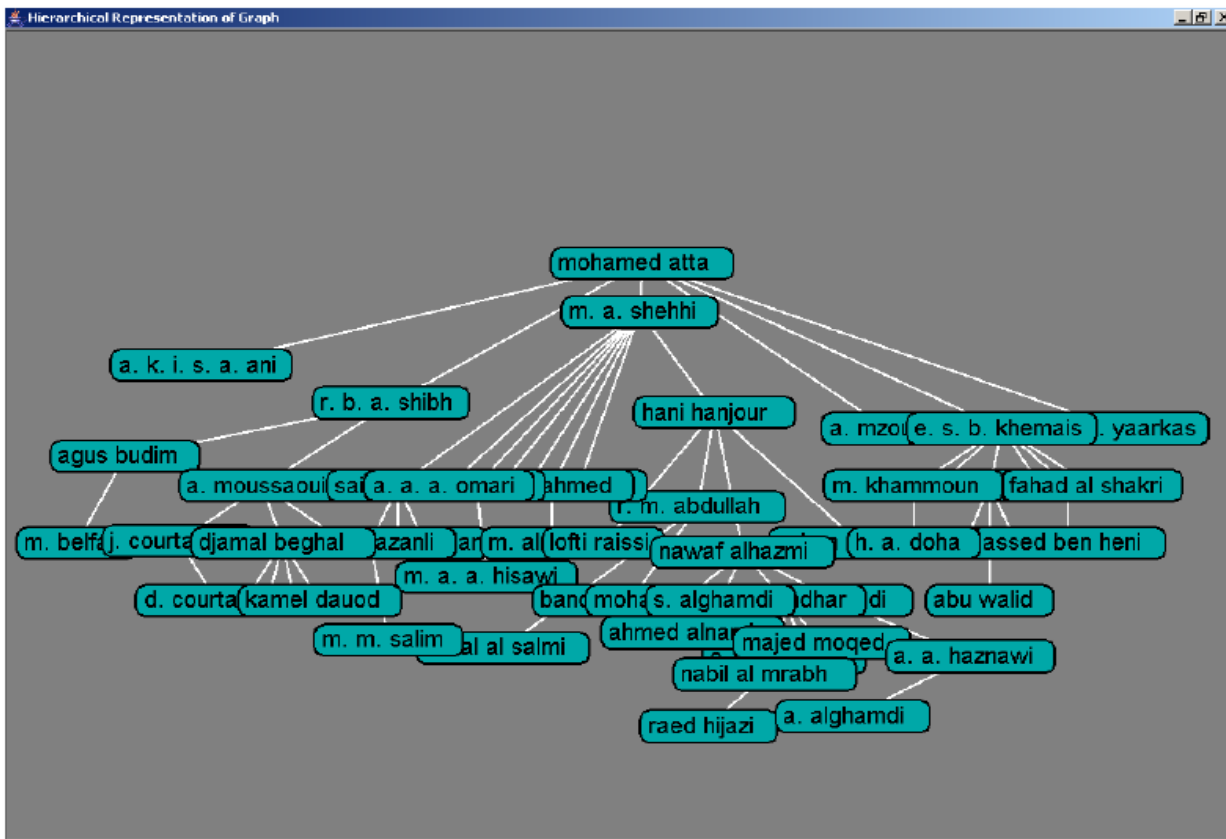
It is fact that Ramzi Yousef began in 1991 to plan bombing attack with USA. Yousef's Uncle Khalid Sheikh Mohammed gave him advice and tips over the phone, and funded him. Yousef entered USA with a Iraqi passport without a U.S. visa on September 1, 1992 [35]. Muhammad Jamal Khalifa, Omer Abdel Rehman, Abdul Rehman Yasin, Abdul Hakim Murad were key companions of Ramzi Yousef. There is an error found in the dataset that Ahmed Yousef is an alias used by Ramzi Yousef, but mistakenly it appears as an independent actor (This is an error). Moreover M. Salameh was an informant and he informed FBI about the plot.

#### **6.4 September 11, 2001 Terrorist Plot**

The September 11, 2001 attacks (often referred to as 9/11—pronounced "nine eleven") consisted of a series of coordinated terrorist suicide attacks upon the United States, predominantly targeting civilians, carried out on Tuesday, September 11, 2001 [36].

That morning, 19 terrorists affiliated with al-Qaeda hijacked four commercial passenger jet airliners. Each team of hijackers included a trained pilot. Two aircraft (United Airlines Flight 175 and American Airlines Flight 11) crashed into the World Trade Center in New York City, one plane into each tower (WTC 1 and WTC 2). Both towers collapsed within two hours, followed by WTC 7 later that day. The pilot of the third team crashed American Airlines Flight 77 into the Pentagon in Arlington County, Virginia. Passengers and members of the flight crew on the fourth aircraft (United Airlines Flight 93) attempted to retake control of their plane from the hijackers; that plane crashed into a field near the town of Shanksville in rural Somerset County, Pennsylvania. As well as the 19 hijackers, a confirmed 2,973 people died and another 24 are missing but presumed dead as a result of these attacks [36].

The renowned Social Network Analyst Valdis Krebs [12] mapped the network of 9/11 (hijackers and their affiliates) as shown in Figure 1. Using the algorithms for detecting hidden hierarchy of non-hierarchical terrorist networks [2], we tested the network of terrorists involved in 9-11 tragic events and results are depicted in Figure 15.



**Figure 15. The hierarchy clearly suggests that Muhammad Atta was the most powerful person (leader) of the plot. While M.A. Shehhi was assisting him, as he is below in the hierarchy. They both were found the key leaders of the plot by 9-11 commission.**

## 7.0 CONCLUSION

In this paper, we presented an overview of an investigative data mining toolkit (iMiner software prototype) which we have developed for undertaking analysis of terrorist networks. In general investigative data mining has been shown to be a promising and potentially powerful area of research. The paper presented interesting patterns gleaned from the data. We discussed three innovative ideas of our research which were already published in [2] and featured in Government Computer News [31]. The mathematical models and algorithms discussed in the paper are implemented in the prototype. The *iMiner* demonstrates key capabilities and concepts of a terrorist network analysis toolkit. Using the toolkit investigating officials can predict overall functionality of the network along with key players. Thus counterterrorism strategy can be designed keeping in the mind that destabilization not only means disconnecting the dots (nodes) but disconnecting those key players from the peripheries by which maximum network could be disrupted. The investigative mining can be used to understand terrorist networks, and we are of the view that investigative data mining tool like iMiner, score over traditional analysis of networks with large volume of data and investigative data mining could reduce the consequent overload on analysts. The results presented in this paper are our findings based on limited exercise in exploring the utility of investigative data mining in analyzing terrorist networks. The may also be used for law enforcement agencies for destabilizing of terrorist networks for capturing the key nodes. The intelligence agencies may also evaluate the efficiency of the networks by capture of a particular node.

Further real-time or near real-time information from multiplicity of databases could have the potential to generate early warning signals of utility in detecting and deterring terrorist attacks. It is necessary, of course, to have ‘experts’ in the loop. This analysis has provided substantive and in-depth analysis of terrorist networks. Furthermore this analysis has provided a richer and deeper understanding and insight into terrorist networks and has provided approaches to destabilize the networks.

In this paper we presented the system architecture of our software prototype and the process by which we harvested data from web and stored in the knowledgebase. The focus of the knowledge base we have developed is the agglomeration of publicly available data and integration of the datasets with the software prototype in order to investigate interested patterns. The flow model of the prototype is also shown and discussed

## **8.0 REFERENCES**

- [1] Memon Nasrullah and Henrik Legind Larsen.: *Practical Approaches for Analysis, Visualization and Destabilizing Terrorist Networks*. In the proceedings of ARES 2006: The First International Conference on Availability, Reliability and Security, Vienna, Austria, IEEE Computer Society, pp. 906-913, (2006)
- [2] Memon Nasrullah and Henrik Legind Larsen: *Practical Algorithms of Destabilizing Terrorist Networks*. In the proceedings of IEEE Intelligence Security Conference, San Diego, Lecture Notes in Computer Science, *Springer-Verlag*, Vol. **3976**: pp. 398-411 (2006)
- [3] Memon Nasrullah and Henrik Legind Larsen: Detecting Terrorist Activity Patterns using Investigative Data Mining Tool. *International Journal of Knowledge and System Sciences*, Vol. **3**, No. 01, pp. 43-52 (2006)
- [4] Memon Nasrullah, Abdul Rasool Qureshi.: Destabilizing Terrorist Networks. In *WSEAS Transactions on Computers*, Issue 11, Vol. **4**, pp.1649-1656 (2005)
- [5] Taskar Ben, Pieter Abbeel, Ming-Fai Wong, and Daphne Koller: *Label and Link Prediction in Relational Data*, in IJCAI Workshop on Learning Statistical Models from Relational Data, (2003)  
[http://kdl.cs.umass.edu/srl2003\\_upload/files/taskar-paper.pdf](http://kdl.cs.umass.edu/srl2003_upload/files/taskar-paper.pdf)
- [6] M. Barlow, J. Galloway, and H. Abbass, *Mining Evolution through Visualization*. In Proceedings of Workshop on Beyond Fitness: Visualization Evolution at the 8th International Conference on the Simulation and Synthesis of Living Systems, (2002)  
<http://www.alife.org/alife8/workshops/15.pdf>
- [7] Q&A with Professor Karen Stephenson, April 18, 2006  
[http://www.elearningpost.com/articles/archives/qa\\_with\\_professor\\_karen\\_stephenson/](http://www.elearningpost.com/articles/archives/qa_with_professor_karen_stephenson/)
- [8] Klerks, P. The Network Paradigm Applied to Criminal Organisations: Theoretical nitpicking or a relevant doctrine for investigators? Recent developments in the Netherlands. *Connections* 24(3): 53-65, (2002)
- [9] Fijnaut Cyrille J.C.F., Frank Bovenkerk, Gerben Bruinsma. Organized Crime in Netherlands. *The Hague Kulwar Law International*, (1998)
- [10] Freeman, L.C. Centrality in Social Networks: I. Conceptual clarification. *Social Networks*, **1**:215-39 (1978)
- [11] Bonacich, P. Factoring and weighting approaches to status scores and clique identification. *Journal of Mathematical Sociology*. **2**:113-120 (1972)
- [12] Krebs, V. E. Mapping Network of Terrorist Cells. *Connections* 24(3): 43-52 (2002)
- [13] Allen, C. *The Dunbar Number as a Limit to Group Sizes* (2004). Retrieved May 31, 2006, from [http://www.lifewithalacrity.com/2004/03/the\\_dunbar\\_numb.html](http://www.lifewithalacrity.com/2004/03/the_dunbar_numb.html)
- [14] Watts, D.J.: *Six Degrees – The Science of a Connected Age*. W.W. Norton & Company, New York, 2003.
- [15] Kreisler, H.: International Relations in Information Age: Conversation with John Arquilla, Professor at University of California, Berkeley 2003. Retrieved on June 01, 2006 from <http://globetrotter.berkeley.edu/people3/Arquilla/arquilla-con0.html>
- [16] Al Qaeda Training Manual. Retrieved on June 01, 2006 from



<http://www.fas.org/irp/world/para/manualpart1.html>

- [17] Walter R. Harper and Douglas H. Harris, The Application of Link Analysis to Police Intelligence,” *Human Factors* 17: 157-164 (1975)
- [18] Carley, Kathleen Ju-Sung Lee, David Krackhardt. Destabilizing Networks, *Connections* 24(3):31-34. (2001)
- [19] McAndrew, D.: *Structural Analysis of Criminal Networks. Social Psychology of Crime: Groups, Teams, and Networks, Offender Profiling Series*, III. L. Allison. Dartmouth, Aldershot (1999)
- [20] Sparrow, M. K.: Application of Network Analysis to Criminal Intelligence: An Assessment of the Prospects. *Social Networks* 13 , 251-274 (1991)
- [21] Sageman, M.: Understanding Terror Networks. Pennsylvania, University of Pennsylvania Press (2004)
- [22] J. Raab and H. B. Milward. Dark Networks as Problems. *Public Administration Research and Theory*, 13(4): 413-439 (2003)
- [23] P. Svenson, C. Mårtensson, C. Carling. *Complex Networks: Models and Dynamics*, FOI-R-1766-SE (2005)
- [24] A. Clauset, M. Newman, C. Moore. Finding Community Structure in Very Large Networks, *Physical Review E* 70, 066111 (2004)
- [25] M. Newman. A Measure of Betweenness Centrality based on Random Walks. *Social Networks* 27: 39-54 (2005)
- [26] Farley, J.D. Breaking Al Qaeda Cells: A Mathematical Analysis of Counterterrorism Operations (A Guide for Risk Assessment and Decision Making). *Studies in Conflict and Terrorism*, 26: 399-411 (2003)
- [27] Latora, V., Massimo Marchiori, How Science of Complex Networks can help in developing Strategy against Terrorism, *Chaos, Solitons and Fractals* 20, 69-75 (2004)
- [28] Bali Night Club Bombing. Retrieved on June 01, 2006 from <http://www.globalsecurity.org/security/ops/bali.htm>
- [29] Karl M. Van Meter, “Terrorists/Liberators: Researching and Dealing with adversary social networks,” *Connections* 24 (3): 66-78 (2001)
- [30] Adnan el Shukrijumah is a Most Wanted Terrorist  
<http://www.ianlivingston.com/archived/2004/shukrijumah.htm>
- [31] Jackson Joab, NSA and Social Networking: TRENDS & TECHNOLOGIES that affect the way government does IT, Government Computer News, 24<sup>th</sup> July 2006. Available online  
[http://www.gcn.com/print/25\\_21/41403-1.html](http://www.gcn.com/print/25_21/41403-1.html)
- [32] i2: Analyst’s Notebook (Accessed on November 14, 2006).  
[http://www.i2.co.uk/Products/Analysts\\_Notebook/default.asp](http://www.i2.co.uk/Products/Analysts_Notebook/default.asp)
- [33] Wheeler, L. Scott , Dirty Bomb Plot underway in US? (Accessed on November 14, 2006)  
[http://www.worldnetdaily.com/news/article.asp?ARTICLE\\_ID=35339](http://www.worldnetdaily.com/news/article.asp?ARTICLE_ID=35339)
- [34] Adnan Gulshair el Shukrijumah (Accessed on November 14, 2006)  
[http://en.wikipedia.org/wiki/Adnan\\_Gulshair\\_el\\_Shukrijumah](http://en.wikipedia.org/wiki/Adnan_Gulshair_el_Shukrijumah)
- [35] Ramzi Yousef (Accessed on November 14, 2006)  
[http://en.wikipedia.org/wiki/Ramzi\\_Yousef](http://en.wikipedia.org/wiki/Ramzi_Yousef)
- [36] September 11, Attacks (Accessed on November 14, 2006)  
[http://en.wikipedia.org/wiki/September\\_11,\\_2001\\_attacks](http://en.wikipedia.org/wiki/September_11,_2001_attacks)





# **Investigative Data Mining Toolkit:** **A Software Prototype for Visualizing, Analyzing and Destabilizing Terrorist Networks**

October 20, 2006



*Nasrullah Memon*





**Aalborg University Denmark**  
Software Intelligence Security Research Center

# Outline

- Introduction
- Background
- Challenges
- Investigative Data Mining Toolkit
- Power Analysis
- Role Analysis
- Cohesion Analysis
- Detection of Hidden Hierarchy
- Case Studies
- Conclusion





# Introduction

- The concern about national security has increased significantly since the terrorist attack on September 11, 2001
  - Intelligence agencies such as the CIA, FBI and NSA are actively collecting and analyzing information to investigate terrorists' activities
  - Academic world has increased the attention paid to analyzing terrorist networks. The analysis of terrorist networks divided into two groups: Data Collectors and Modelers
- 
- 





# Background

- After tragic attacks by kidnapped airlines on New York and Washington in **September 2001**, the interest for Al Qaeda in public and media rose immediately.
  - Experts and Analysts all over the world started to offer various explanations of Al Qaeda's origins, membership recruitment, modes of operation, as well as **possible ways of disruption**.
  - Al Qaeda is "a net that contains independent intelligence", that it "functions as swarm", that it "gathers from nowhere and disappears after action", that it is "an **ad hoc network**", "an **atypical organization**".
  - America's intelligence community stands at a **critical crossroads**, there is need how best to improve the collection and analysis of critical foreign intelligence as America fights an increasingly **dangerous** international war on terror.
  - Current terrorist threat is not organized with **conventional lines** of authority.
- 
- 



# Background

- Instead they are organized as **loose networks** and so belong to an analytically **distinct category**.
- Al Qaeda has evolved from a **centrally directed organization** into a worldwide **franchiser** of terrorist attacks
- Al Qaeda does appear to have become increasingly **decentralized** after the arrest of key Al Qaeda leaders and losing the safe haven i.e. Afghanistan.
- Al Qaeda convened a **strategic summit** in northern Iran in November 2002, at which it was decided that it could no longer operate as **hierarchy**, but instead would have to decentralize.
- Looking to the facts and figures, we propose **mathematical models and practical algorithms** for Destabilizing Terrorist Networks.
- In our research we have also constructed a **knowledgebase** of the terrorist events occurred in the past.

# Challenges

- **Data Collection** is difficult for any network analysis because it is hard to create a complete network
  - It is specially difficult to gain information on terrorist networks
    - **Terrorist organizations** do not provide information on their members, and **government** rarely allows researchers to use their data
  - A **number of academic researchers** focus primarily on data collection on terrorist organizations, analyzing the information through description and straightforward modeling
    - E.g. Valdis Krebs, Jose A. Rodriguez, Marc Sagman
      - Despite **many strengths**, there are a **few drawbacks**
      - By dealing with open sources, these authors are limited in acquiring data. With open sources, if author does not have information on terrorists, he or she assumes they do not exist. If researcher could not find an al- Qaeda operative in US, he could assume, no al-Qaeda network in US



# Challenges

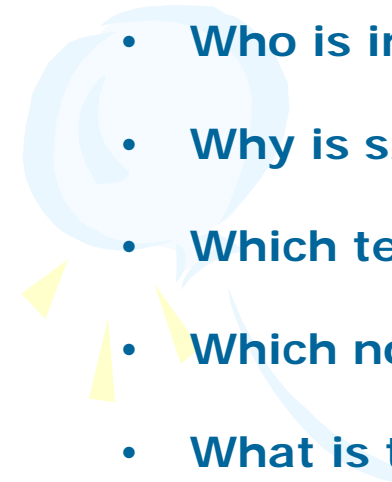
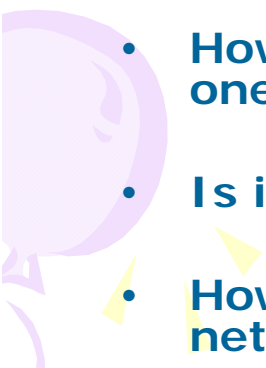
- The **common problem** for the modelers (in the field of analysis of terrorist networks) is the **issue of data**
- Any academic work is **only as good as the data**, no matter the type of advanced methods used
- Modelers often do not have the **best data**, as they have not collected individual biographies (like sageman) and do not have access to classified data
- Many of the models are created data-free or without complete data
- Modelers often do not have a foundation in terrorist studies nor do they always work with top counter-terrorism experts
  - Without the help of **counterterrorism experts**, it is difficult to turn the numbers and the graphic models into interpretable results that make sense in the context of vast literature on terrorism



# Goals

Intelligence and law enforcement agencies are often interested in finding structural properties of terrorist networks.



**This study aims to answer the following questions:**

- How **Investigative Data Mining** will be a useful tool for Law Enforcement and Intelligence agencies in War against Terrorism?
  - Who is important in a network?
  - Why is s/he **important**?
  - Which terrorist is **highly/less** connected?
  - Which nodes (terrorists) are **key players**?
  - What is the **effeciency** of a network?
  - How much the effeciency of the network is **reduced** by eradicating one (or some) of the key player (s)?
  - Is it possible to **construct hierarchy** of non-hierarchical networks?
  - How can the **law enforcement** use (often incomplete and faulty) network data to **disrupt and destabilize** terrorist networks?
- 
- 





# Investigative Data Mining

- **Investigative Data Mining (IDM)** offers the ability to firstly **map** a covert cell, and to secondly **measure** the specific structural and interactional criteria of such a cell.
  - This framework aims to **connect the dots** between individuals and “**map and measure complex and covert terrorist networks**”.
  - The IDM focuses on uncovering the patterning of people's interaction, and correctly interpreting these networks assists “in **predicting behaviour and decision-making within the network**”.
  - The IDM also endows the analyst the ability to measure the **efficiency** of the cell as a whole, and also the **level of activity, ability to access others, and the level of control over a network** each individual possesses.
  - The measurement of these criteria allows specific counter-terrorism applications to be drawn, and assists in the assessment of the most **effective methods of disrupting and destabilizing** a terrorist cell.
- 
- 



# Main Contributions

- **Cohesion Analysis** (Discover tightly coupled nodes in various subgroups using Cliques, N-Cliques, K-Cores, K-Clubs and K-Plexes models from graph theory and subgroup Detection using new algorithms)
  - Memon N., HL Larsen (2006) "Detecting Terrorist Activity Patterns using Investigative Data Mining Tool, *International Journal of Knowledge and System Sciences*, Vol. 3(1), 43-52
  - Memon N., HL Larsen (2006) "Structural Analysis and Mathematical Methods for Counterterrorism, In Proc ADMA 2006, LNAI 4093, pp.1073-84.
- **Role Analysis** (Discover role positions within a network, for example, gatekeepers, leaders and followers. Found efficiency of the network and succeeded in discovering how efficiency of the network is decreasing while removal or capture of a terrorist)
  - Memon Nasrullah, HL Larsen, *Practical Algorithms for Destabilizing Terrorist Networks*, In the proceedings of IEEE Intelligence and Security Conference (ISI 2006) San Diego, LNCS 3975, pp. 389-400.
- **Power Analysis** (Discover most powerful nodes in the network and detect commanding structure in a network)
  - Memon Nasrullah, HL Larsen, *Practical Algorithms for Destabilizing Terrorist Networks*, In the proceedings of IEEE Intelligence and Security Conference (ISI 2006) San Diego, LNCS 3975, pp. 389-400.
  - Memon N. et al., *Detecting Hidden Hierarchy of Terrorist Networks*, In post conference proc DCMMC-2006, Rayburn House US Capitol, Washington, 28-29 September 2006, Springer LNCS.



# Centrality Measures

## Calculates

## Measures

Degree

The number of direct connections to other nodes

Connections to others;  
network activity; power

Closeness

Inverse of the sum of the shortest paths to all other nodes in the network

Members key to network;  
communication; reach;  
reachability

Betweenness

Proportion of times a node is on the shortest path b/w other pairs of nodes

Information control; role as  
intermediary, gatekeeper

EV Centrality

Centrality of centrality  
member must be connected to at least  $g-k$  other members

Overall importance to the  
network; how close one node  
to others who are also close to  
others



# Dependence Centrality (DC)

The dependence centrality of a node is defined as how much the node is dependent on any other node in the network.

This measure shows that how much one node is dependent on the another node. We can also say that how much one node is useful to another node in order to communicate with other nodes of the network.

Mathematically it can be written as:

$$DC_{mn} = \sum_{m \neq p, p \in G} \frac{d_{mn}}{N_p} + \Omega$$

Where **m** is the root node which depends on **n** by  $DC_{mn}$  centrality and  $N_p$  actually is the Number of geodesic paths coming from **m** to **p** through **n**, and  $d_{mn}$  is geodesic distance from **m** to **n**. The  $\Omega$  is taken **1** if graph is connected and **0** in case it is disconnected. In this paper we take  $\Omega$  as 1, because we consider that graph is connected. The first part of the formula tells us that:

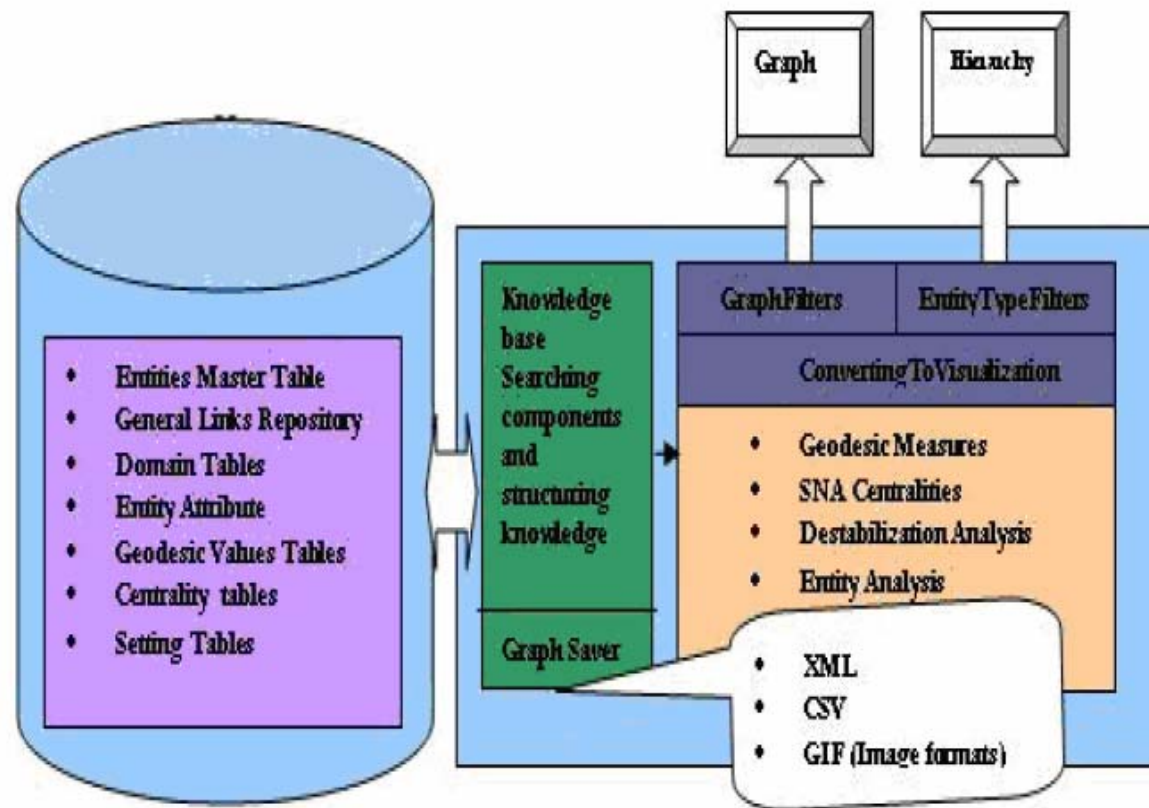
How many times **m** uses **n** to communicate other node **p** of the network?

# Detection of Hidden Hierarchy

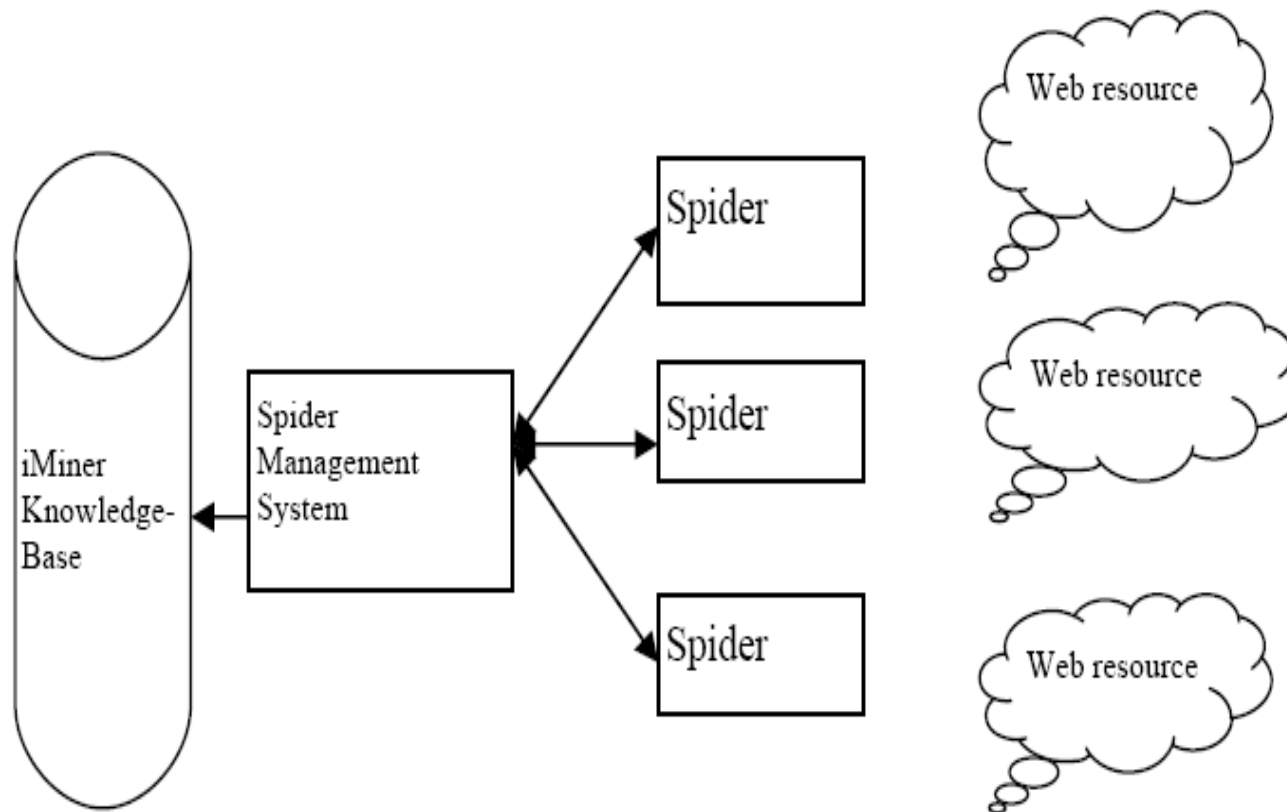
- Using undirected graph, we first convert it into directed graph using degree centrality and Eigen Vector Centrality. For Example, if **degree centrality** of one node is **higher** than other, then simply the **directed link** is originated from that node and point towards other. If they are equivalent in terms of degree, the link will originate from the node with higher **EigenVector centrality**. If Eigen Vector centrality values for **both nodes are equal**, then we ignore the link.
- Then we identify the parents and children pairs. For example, if we have **two nodes**, which are competing for being **parent of a node**, then we have to identify its correct parent. The correct parent will be the one which is connected with **maximum neighbours**. This represents the fact that the true leader, with respect to a node, which is more influential on its neighbourhood.
- Then we identify **hierarchical relationship** among the parents of a node
- At the last step, we detect the **parent of the node** (among the possible parents) by using **dependence centrality**.
- When we identify parents, in such a way we traverse all the nodes. Then a tree structure is obtained, which we call **hierarchical chart**.



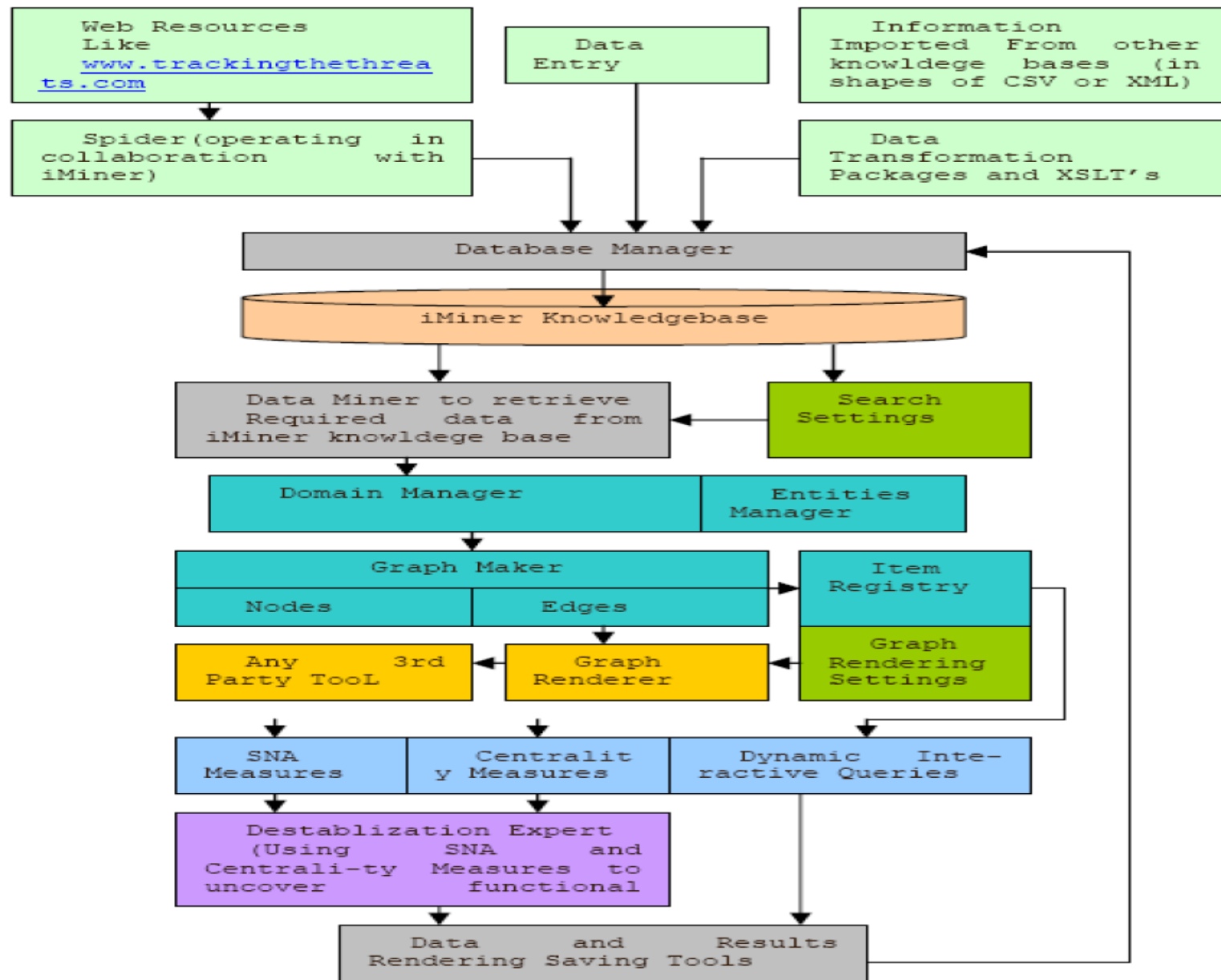
# System Architecture for IDM Toolkit (iMiner)



# Construction of Knowledgebase through Spidering



# Dataflow model





# Role Analysis



# Efficiency of a Network

$$E(G) = \frac{\sum_{i \neq j \in G} \varepsilon_{ij}}{N(N-1)} = \frac{1}{N(N-1)} \sum_{i \neq j \in G} \frac{1}{d_{ij}}$$

The network efficiency  $E(G)$  is a **measure to quantify how efficiently the nodes of the network exchange information.**

To define efficiency of  $G$  first we calculate the shortest path lengths  $\{d_{ij}\}$  between two generic points  $i$  and  $j$ . Suppose that every node sends information along the network, through its edges. The efficiency  $\varepsilon_{ij}$  in the communication between vertex  $i$  and  $j$  is inversely proportional to the shortest distance:  $\varepsilon_{ij} = 1/d_{ij}$

□  $i, j$  when there is no path in the graph between  $i$ , and  $j$ .





# Importance of Node in a Network

The main idea is to use as a measure of the centrality of a node  $i$  the drop in the network efficiency caused by deactivation of the node. The importance  $I(\text{node}_i)$  of the  $i^{\text{th}}$  node of the graph  $G$  is therefore:

$$I(\text{node}_i) \equiv \Delta E = E(G) - E(G - \text{node}_i), i = 1, \dots, N,$$

Where  $G - \text{node}_i$  indicates the network obtained by deactivating node  $i$  in the graph  $G$ . **The most important nodes, i.e. the critical nodes are the ones causing the highest  $\Delta E$ .**

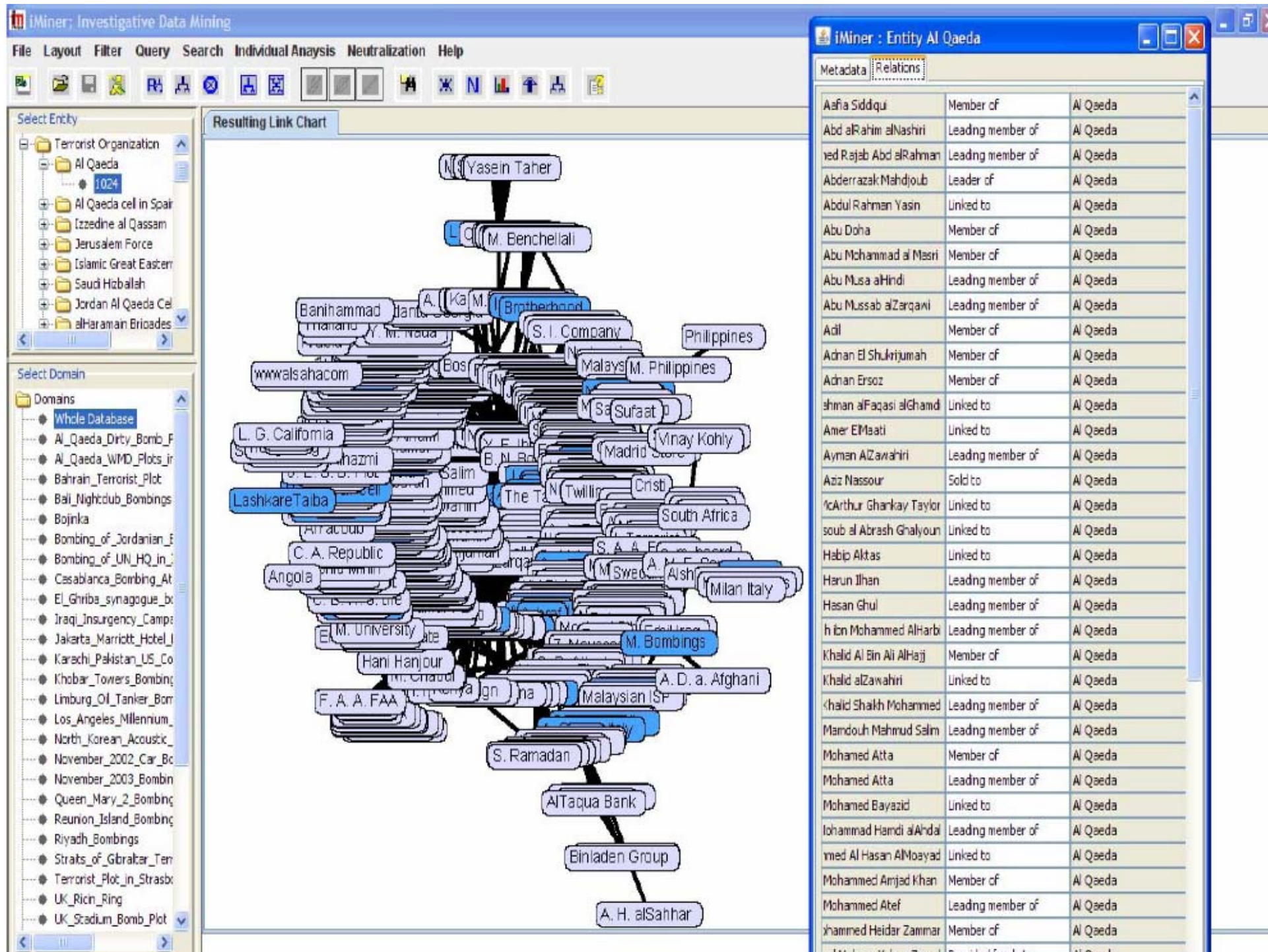


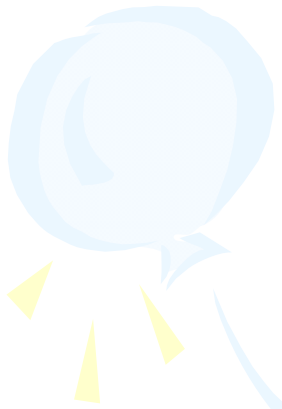
# Position Role Index

- The PRI is proposed measure which highlights a clear distinction between followers and brokers (It is fact that sometime leaders may act as brokers). It depends on the basic definition of efficiency as discussed in last two slides.

$$PRI = E(\dot{G}) - E(\dot{G} - \text{node}_i), i = 1, \dots, N,$$

- It is crystal clear fact that **efficiency of a network** in presence of **followers** is **low** as compared to their absence in the network. This is because they are usually **less connected nodes** and their presence increases the number of low connected nodes in a network, thus **decreasing its efficiency**.
- If we plot the values on the graph, the nodes which are plotted below x-axis are followers, whereas the nodes higher than remaining nodes with higher values on positive y axis are the gatekeepers. While the nodes which are on the x-axis usually central nodes, which can easily bear the loss of any node. The leaders tend to hide on x-axis there.





# **CASE STUDIES**

CASE: 01

## Bali Bombing

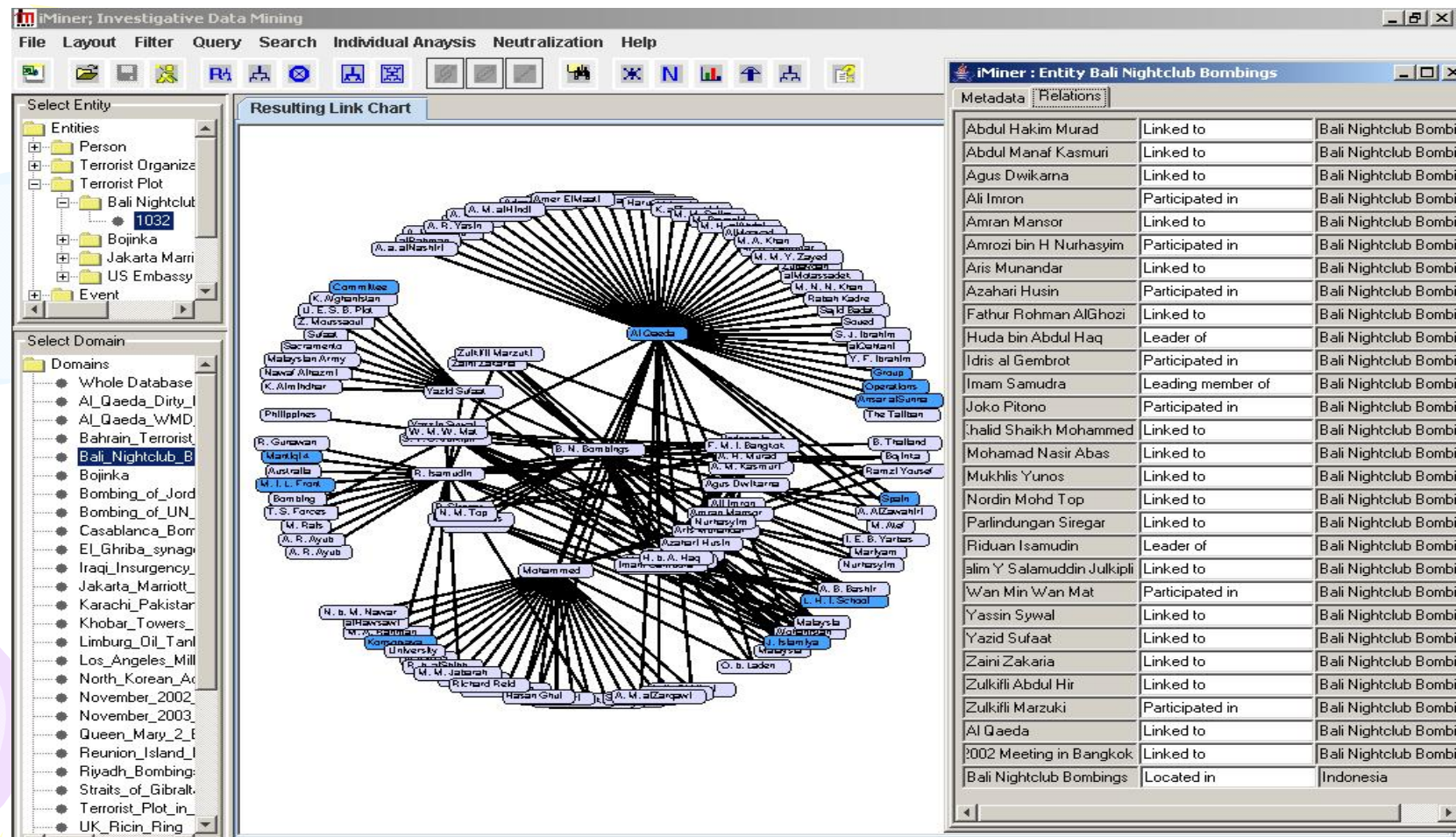




The screenshot displays the iMiner software interface, which is used for investigating data. The interface includes a menu bar with 'File', 'Layout', 'Filter', and 'Query'. Below the menu bar, there are two main windows: 'Select Entity' and 'Select Domain'.

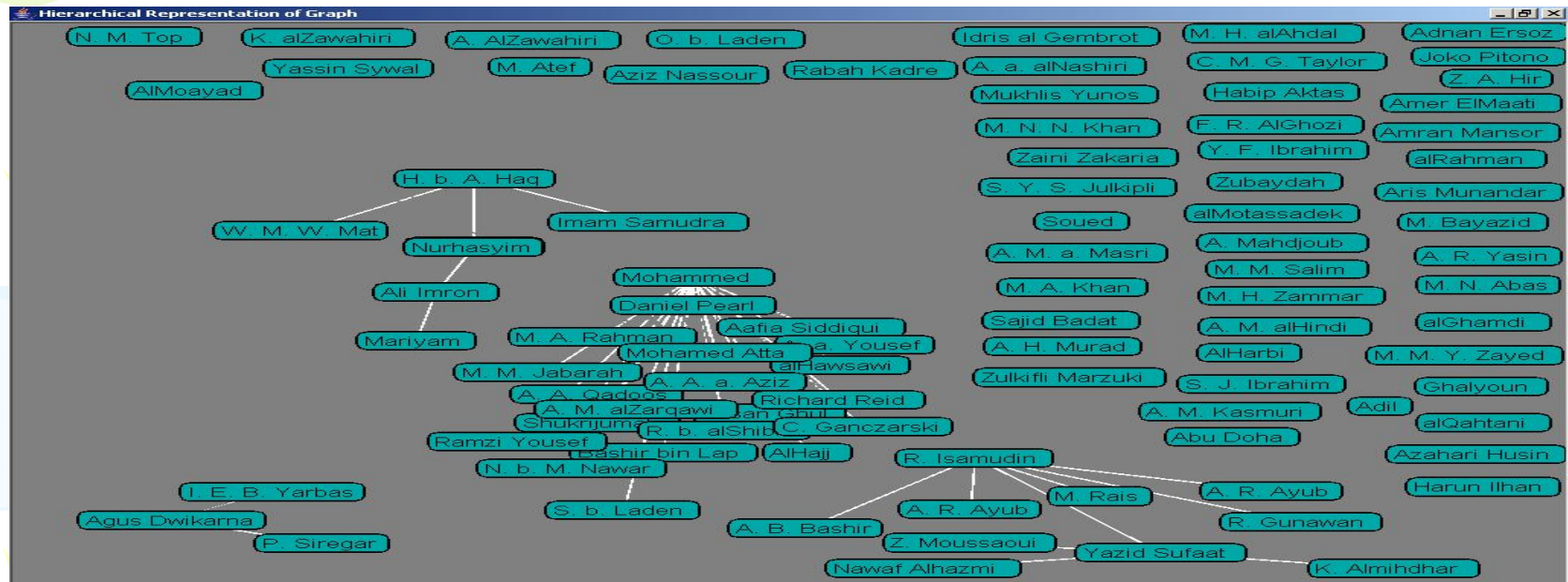
The 'Select Entity' window shows a tree structure of entities. The 'Entities' folder is expanded, revealing sub-folders: 'Person', 'Terrorist Organize', 'Terrorist Plot', 'Bali Nightclub', 'Bojinka', 'Jakarta Marriott', 'US Embassy', and 'Event'. The 'Bali Nightclub' folder is selected, and the number '1032' is displayed next to it.

The 'Select Domain' window shows a list of domains. The 'Domains' folder is expanded, revealing a list of domains: 'Whole Database', 'Al\_Qaeda\_Dirty\_I', 'Al\_Qaeda\_WMD', 'Bahrain\_Terrorist', 'Bali\_Nightclub\_B', 'Bojinka', 'Bombing\_of\_Jord', 'Bombing\_of\_UN', 'Casablanca\_Bom', 'El\_Ghriba\_synag', 'Iraqi\_Insurgency', 'Jakarta\_Marriott', 'Karachi\_Pakistan', 'Khobar\_Towers', 'Limburg\_Oil\_Tan', 'Los\_Angeles\_Mill', 'North\_Korean\_Ac', 'November\_2002', 'November\_2003', 'Queen\_Mary\_2\_E', 'Reunion\_Island\_I', 'Riyadh\_Bombing', 'Straits\_of\_Gibraltar', 'Terrorist\_Plot\_in', and 'UK\_Ricin\_Ring'. The 'Bali\_Nightclub\_B' domain is selected.





# Hierarchy constructed by iMiner for Bali Bombing attack



This hierarchy has some unconnected nodes, where as you can find a hint of patterns some time. The **H. B. A. Haq** and its descendants form a group (This cluster was acted as **executive cluster**), while the cluster **Khalid Shaikh Mohammed** and his affiliates was well known as strategic cluster, whereas **R. Isamudin** (known as Hambali) and his associates cluster known as tactical/logistic cluster .

The accuracy of the algorithms can be determined by the fact that all of **H. B. A. Haq**, **Khalid Shaikh Mohammed** and **R. Isamudin** were key players in the reality. **H. B. A. Haq** was termed as potential leader while **Khalid Sheikh Mohammed** was the key conspirator.

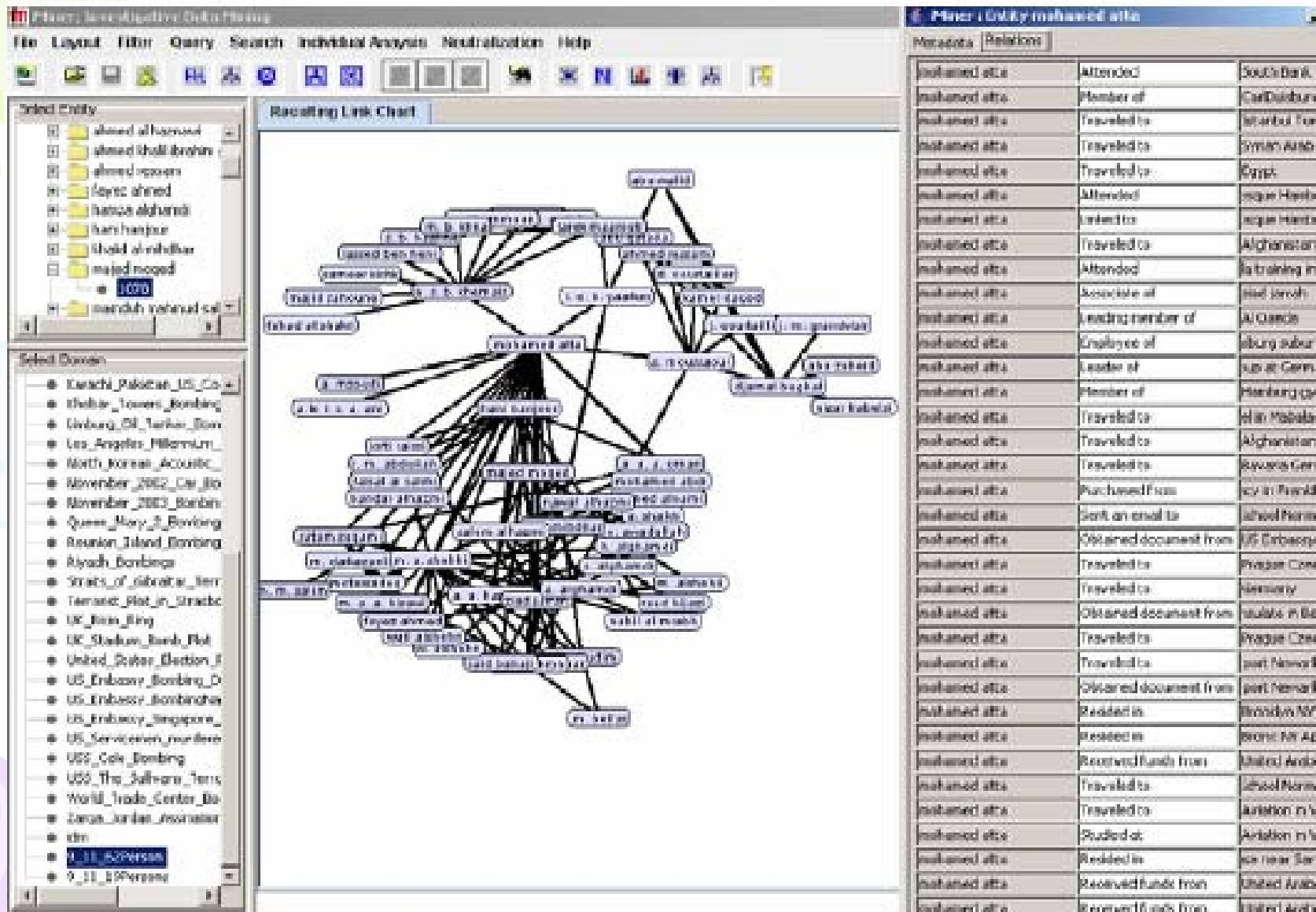
## Effect of Deactivation of a Node of Terrorist Network of Bali Bombing 2002 Attacks

Removed Node	$E(G - \text{Node}_i)$	$\Delta E/E$	K	PRI
K. S. Mohammed	0.306333	0.19844	27	0.1854077
Riduan Isamudin	0.332879	0.12900	23	0.1148397
Yazid Sufaat	0.334258	0.12537	12	0.1111511
Wan Min	0.370011	0.03182	11	0.0160787
Huda Bin A. Haq	0.374484	0.02012	12	0.0041839
Osama Bin Laden	0.374484	0.01868	3	0.0027257

## CASE 02: 9-11 Terrorist Attack

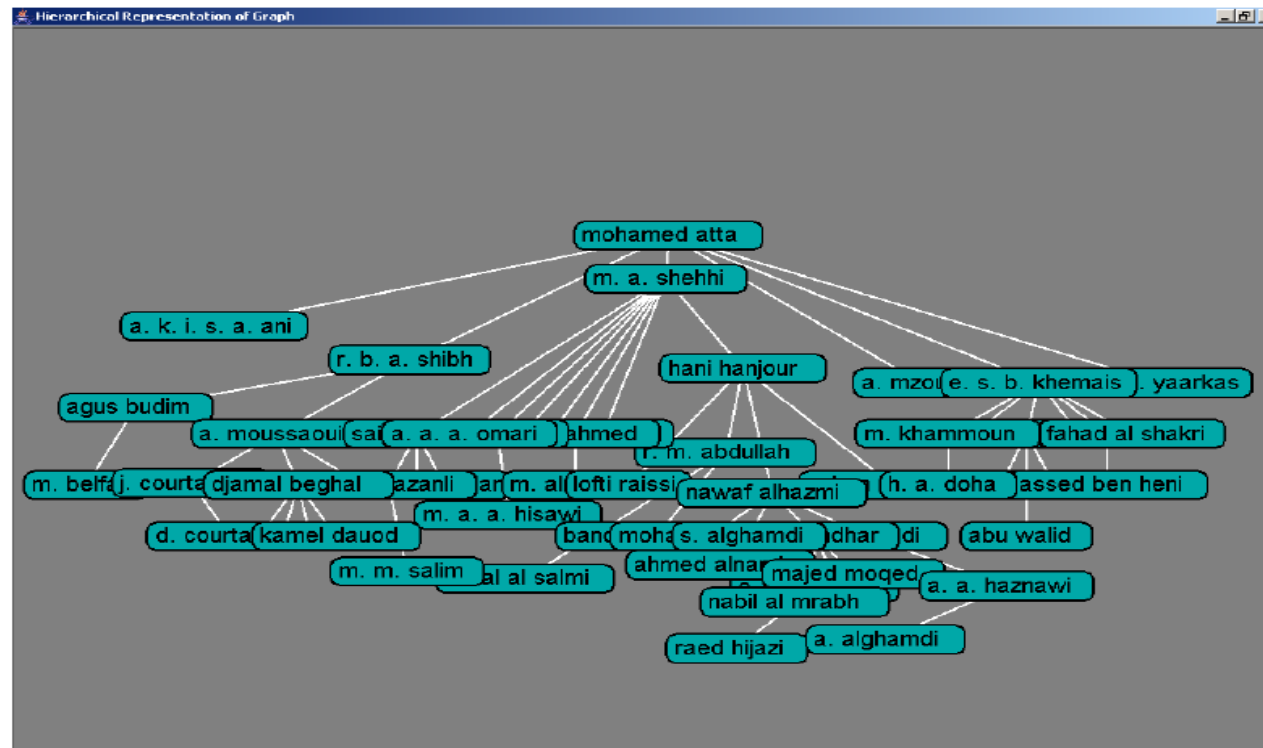


# 9-11 terrorists and affiliates



**The dataset originally designed by Valdis Krebs, but re-constructed in iMiner Using metadata of every terrorist and the event**

# Hierarchical chart for 9-11 hijackers and their affiliates

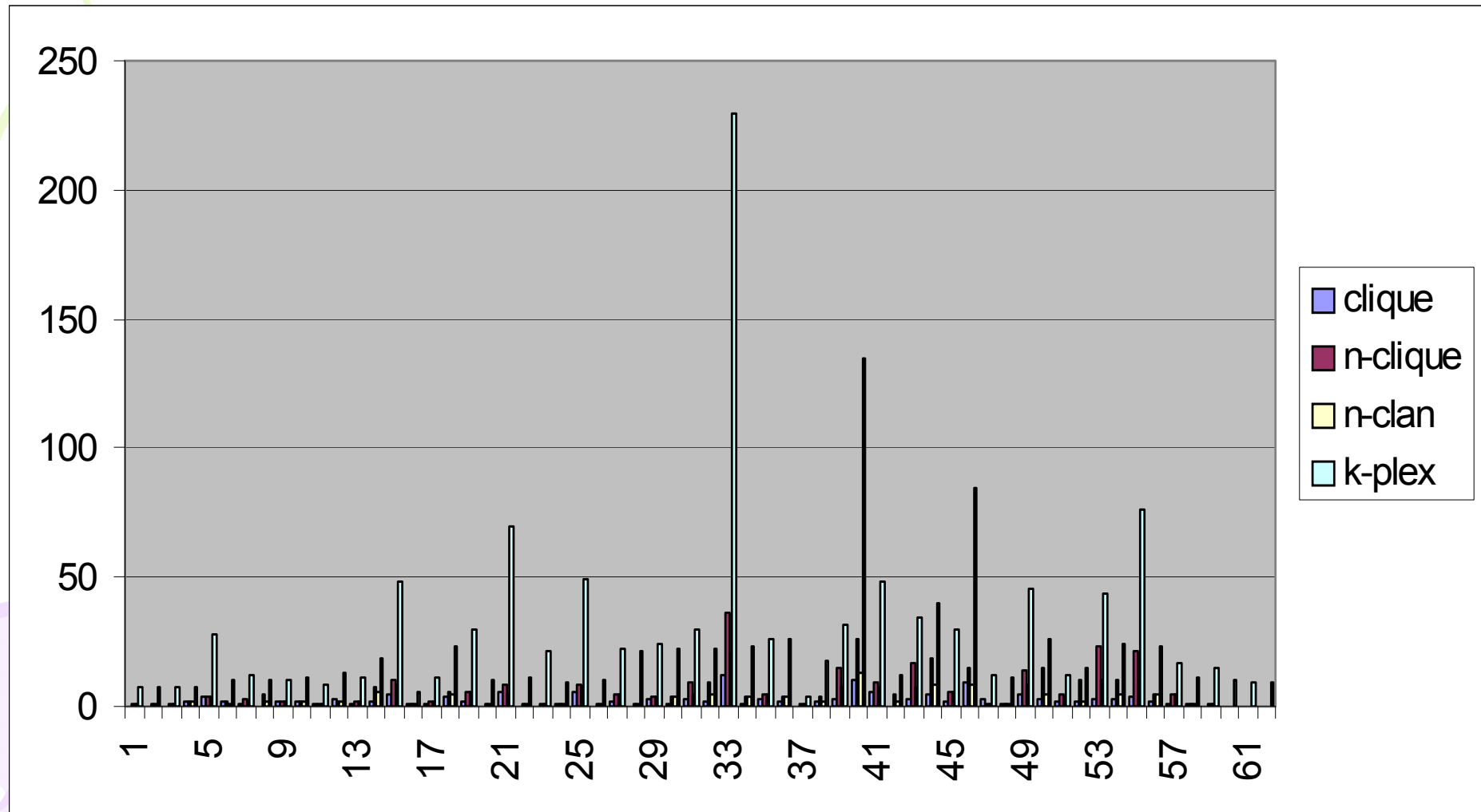


- The hierarchy clearly suggests that **Muhammad Atta** was the **key leader** of the plot. While **Marwan Al Shehri** was assisting him as he is below in the hierarchy. They both were suggested as potential leaders in 9/11 attack and led their respective groups. They were also both members of Hamburg Cell.




## Cohesion Analysis Results of this dataset

# Structural Analysis for 9-11 hijackers and their affiliates

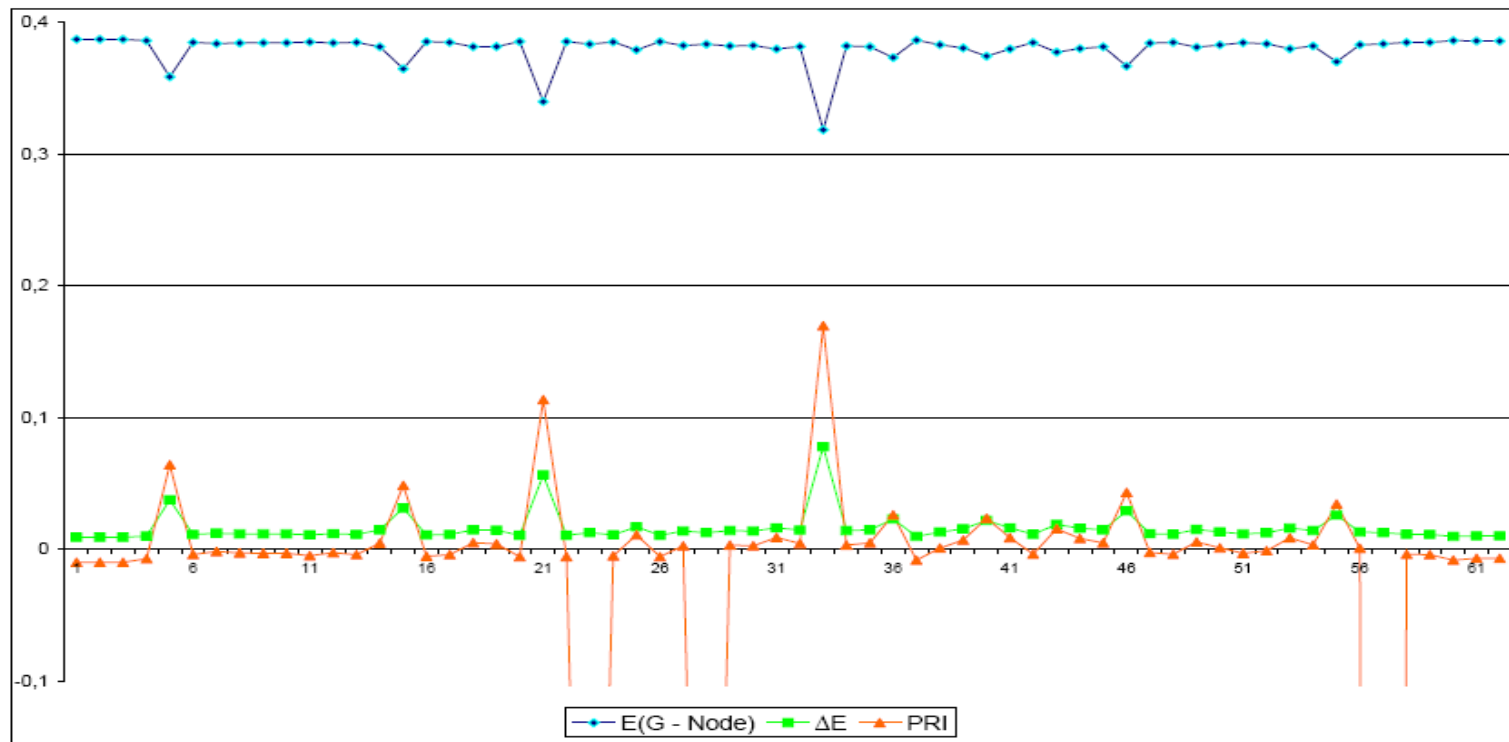






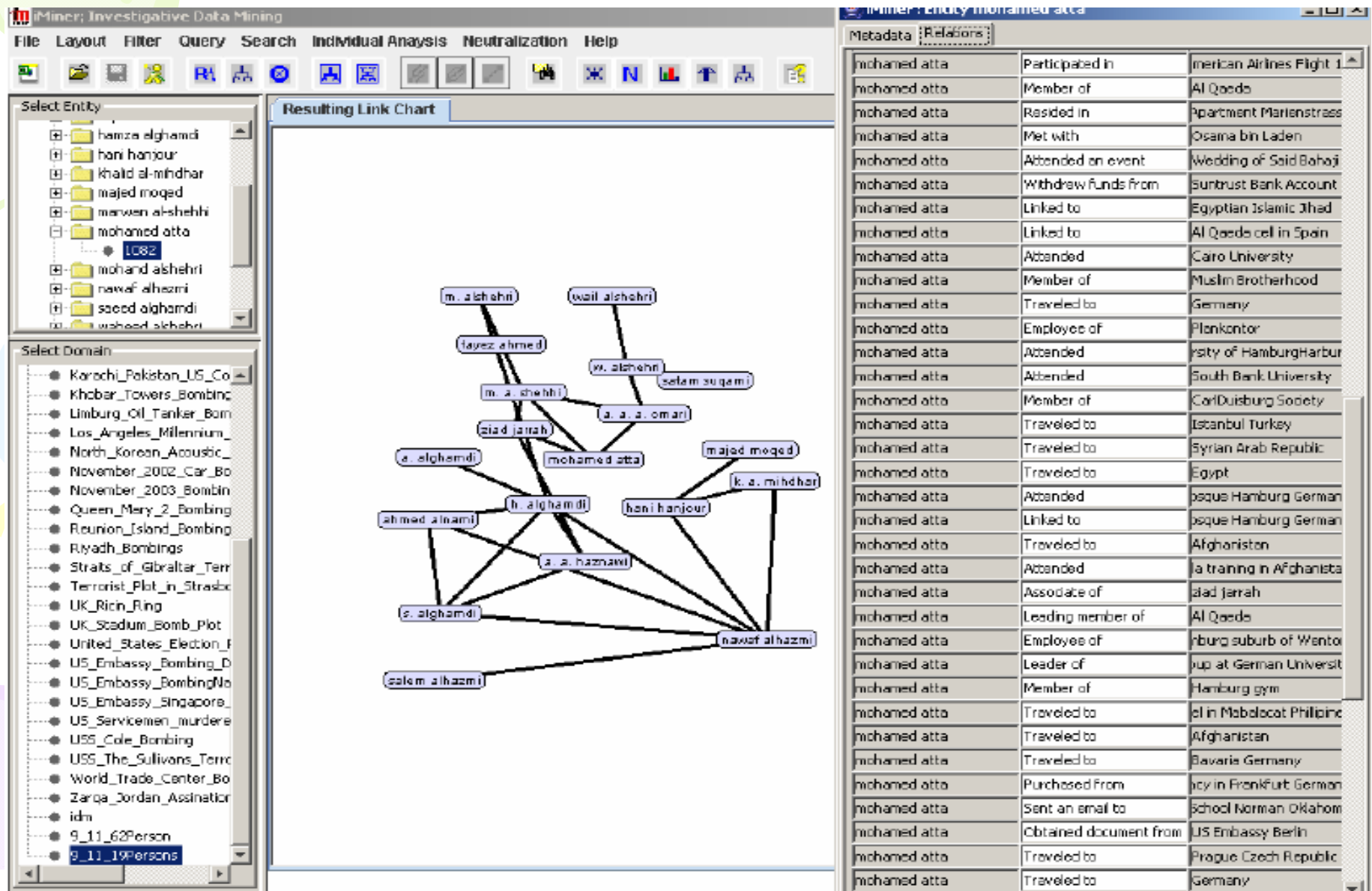
# Role Analysis And Results of the dataset

## Key Players and Important Nodes in 9-11 network (Hijackers and their affiliates)

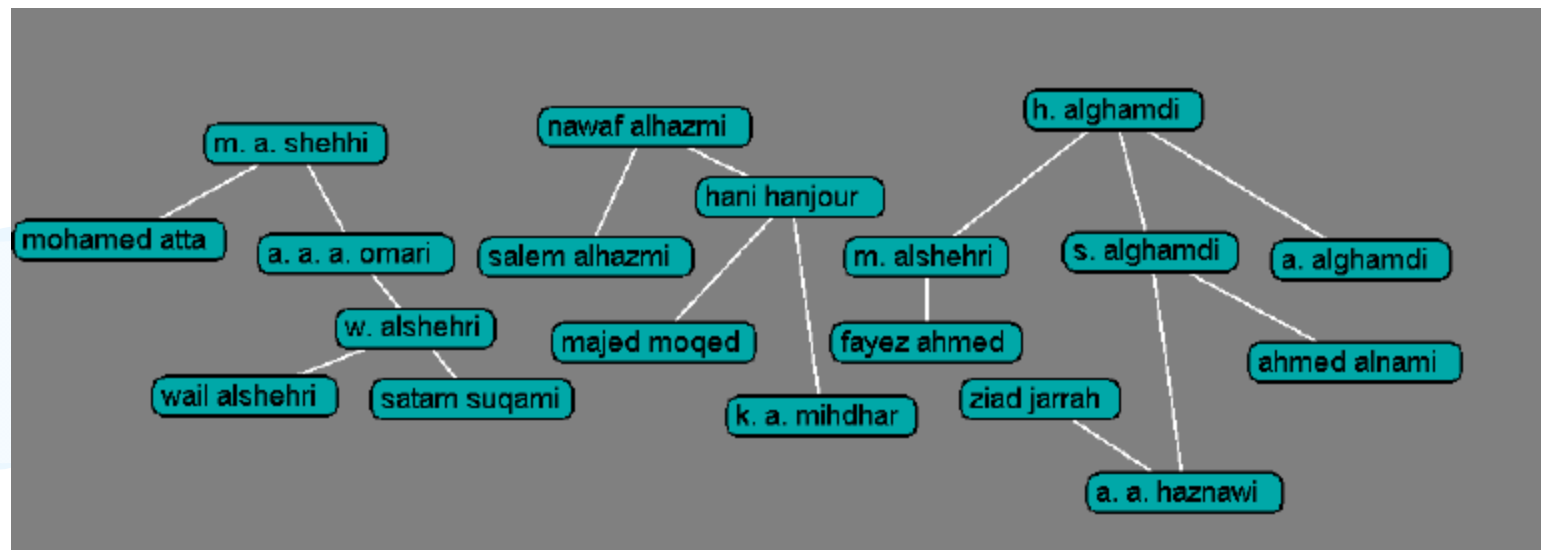


The **efficiency** of the original network is  $E(G) = 0.395$ . The removed node is shown on x-axis; the efficiency of the graph once the node is removed is reported as  $E(G - \text{Node}_i)$ , while the importance of the node (drop of efficiency) is shown as  $\Delta E$ . While position role index shown as PRI of the removed node. The **results prove important aspects of the network** and confirmed that **Mohammed Atta (node # 33)** was the **ring leader**.

# 9-11 Hijackers Terrorist Network



# Hierarchical Chart for 9-11 Hijackers Network



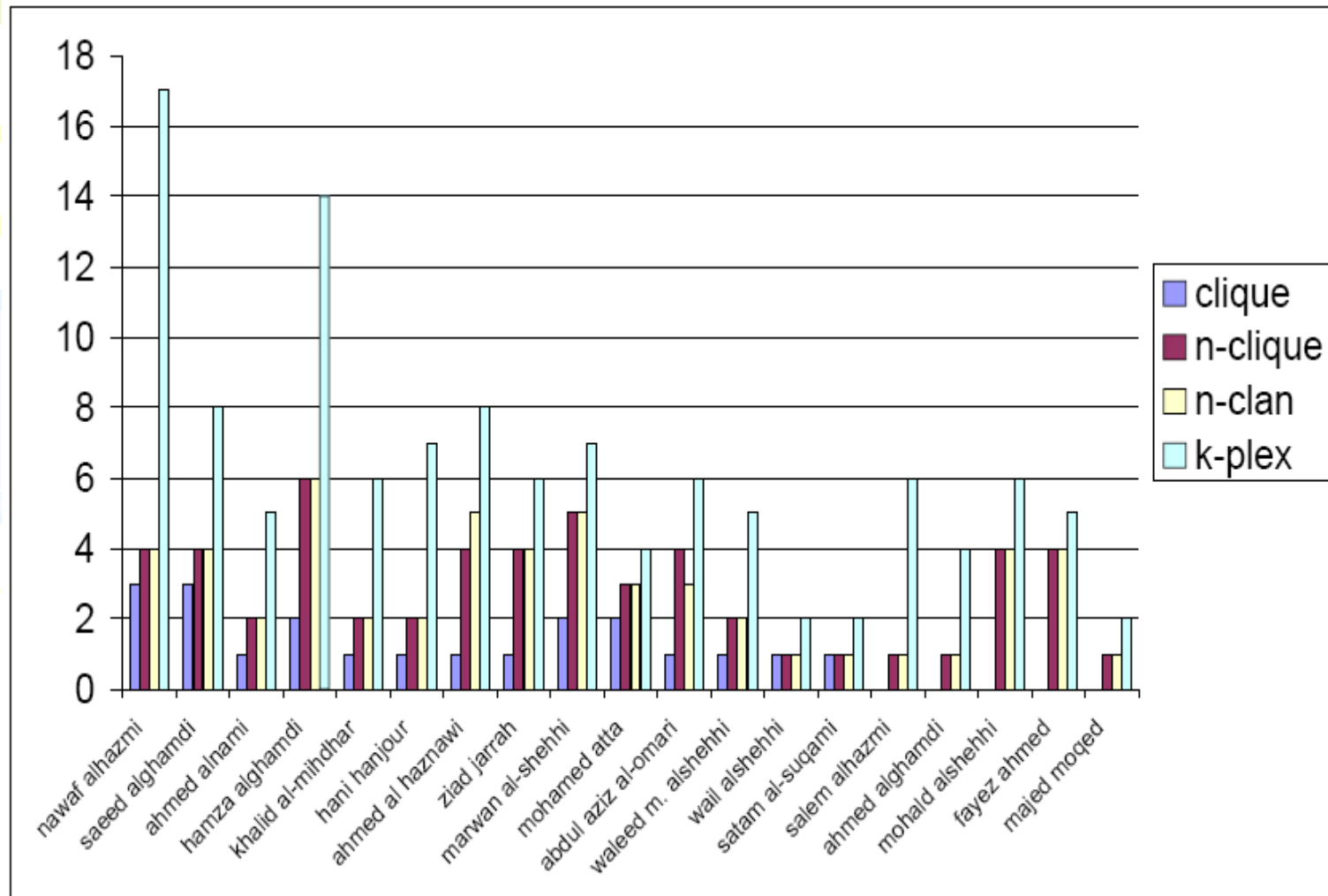
**American Flight 11:** Mohamed Atta, W. Alshehri, Wail Alshehri, a.a. Omari and Satam Suqami

**American Flight 77:** Nawaf Alhazmi, Hani Hanjour, Salem Alhazmi, Majed Moqed and K. Midhar

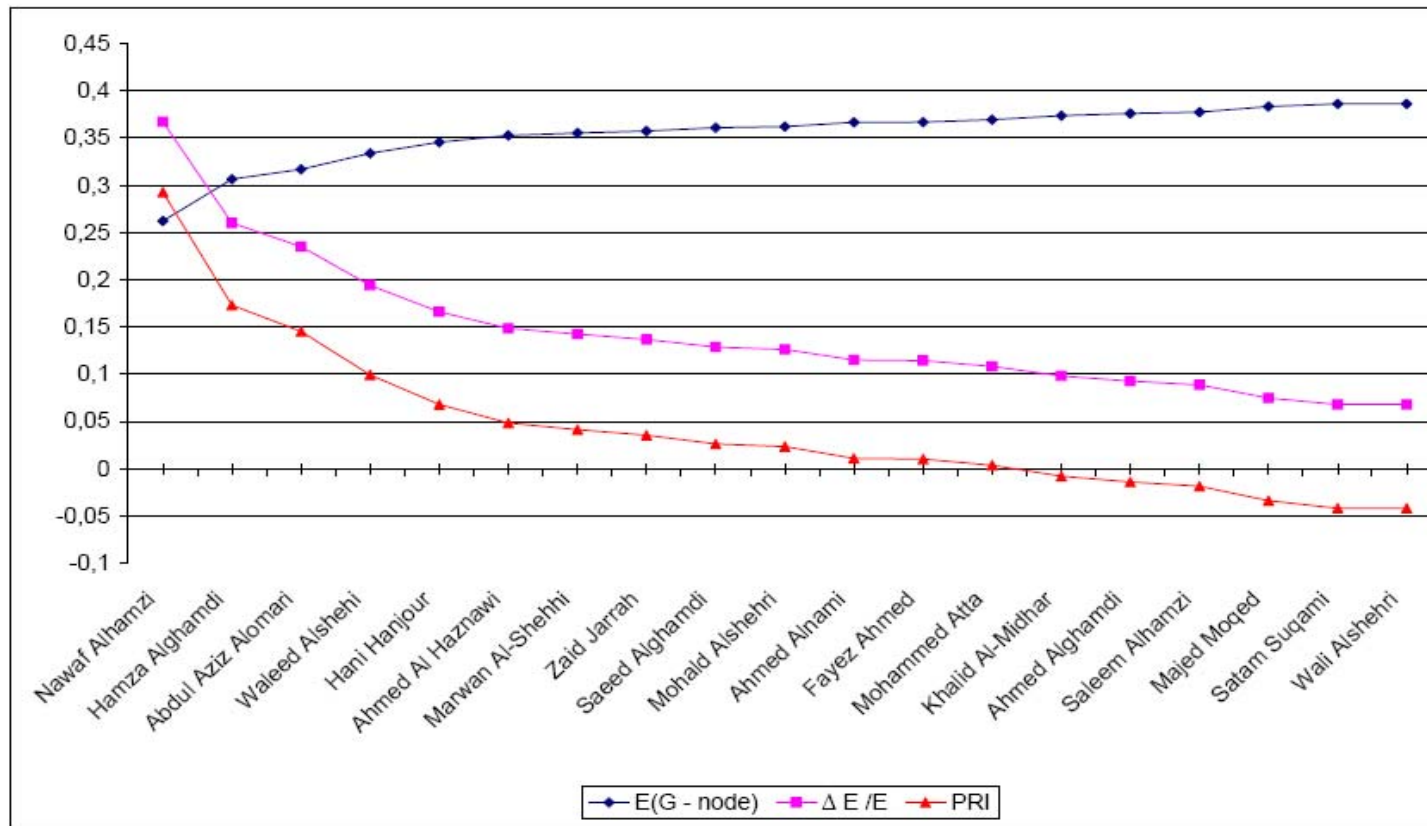
**UA Flight 175:** Marwan Alshehi, Fayz Ahmed, M. Alshehi, Hamza alghamdi and Ahmed alghamdi

**UA Flight 93:** Zaid Jarrah, Ahmed alhaznawi, Saeed Alghmdid and Ahmed alnami

# Structural Analysis for 9-11 hijackers Network



# Key Players and Important Nodes in 9-11 network (Hijackers Network)



The efficiency of the original network  $E(G) = 0.414$ . The removed node is shown on x-axis, the efficiency of the graph once the node is removed is shown as  $E(G - \text{node})$ ; while the relative drop of efficiency is shown as  $\Delta E / E$ . The newly introduced measure position role index is shown as PRI.

CASE 03:

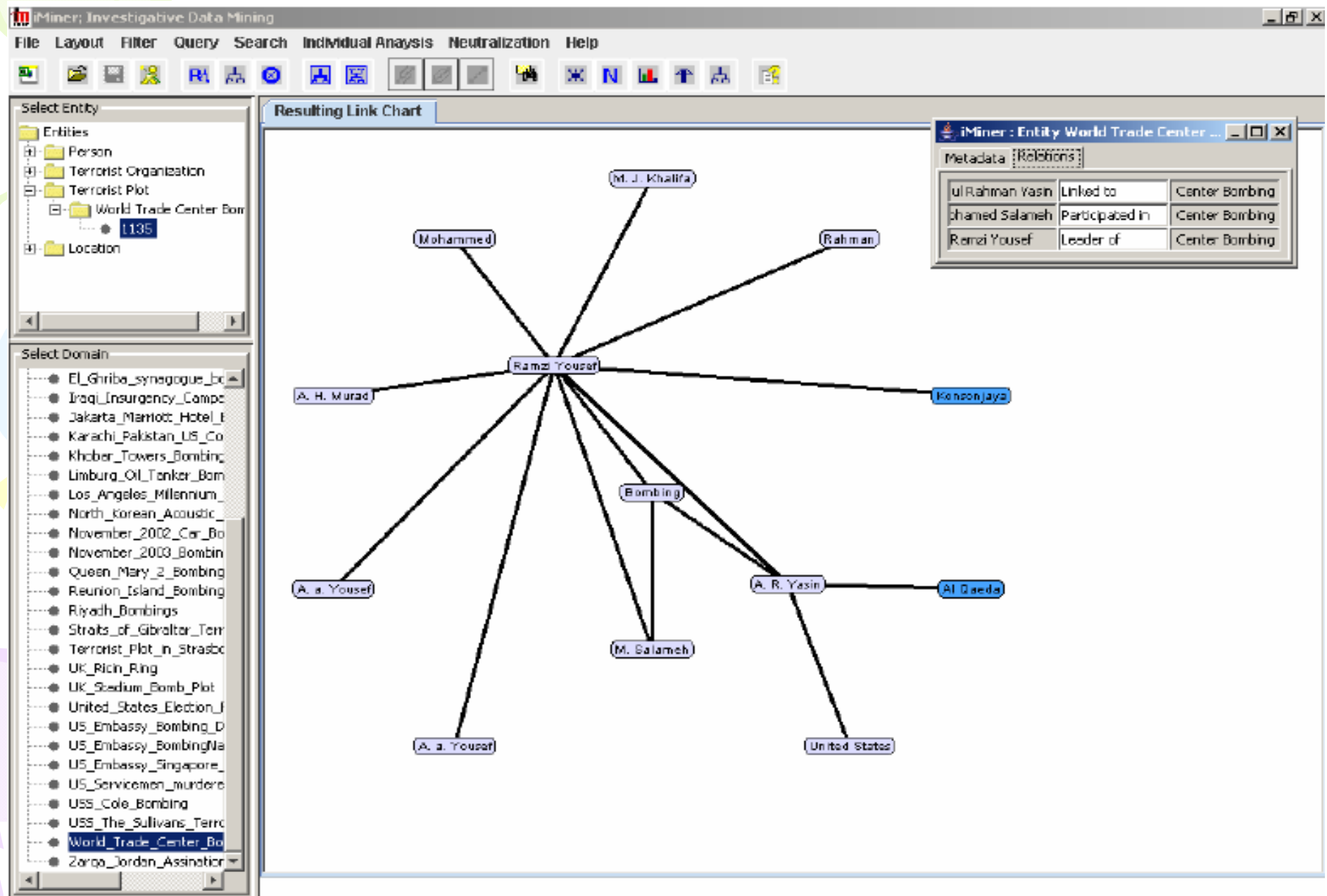
## 1st World Trade Center 1993



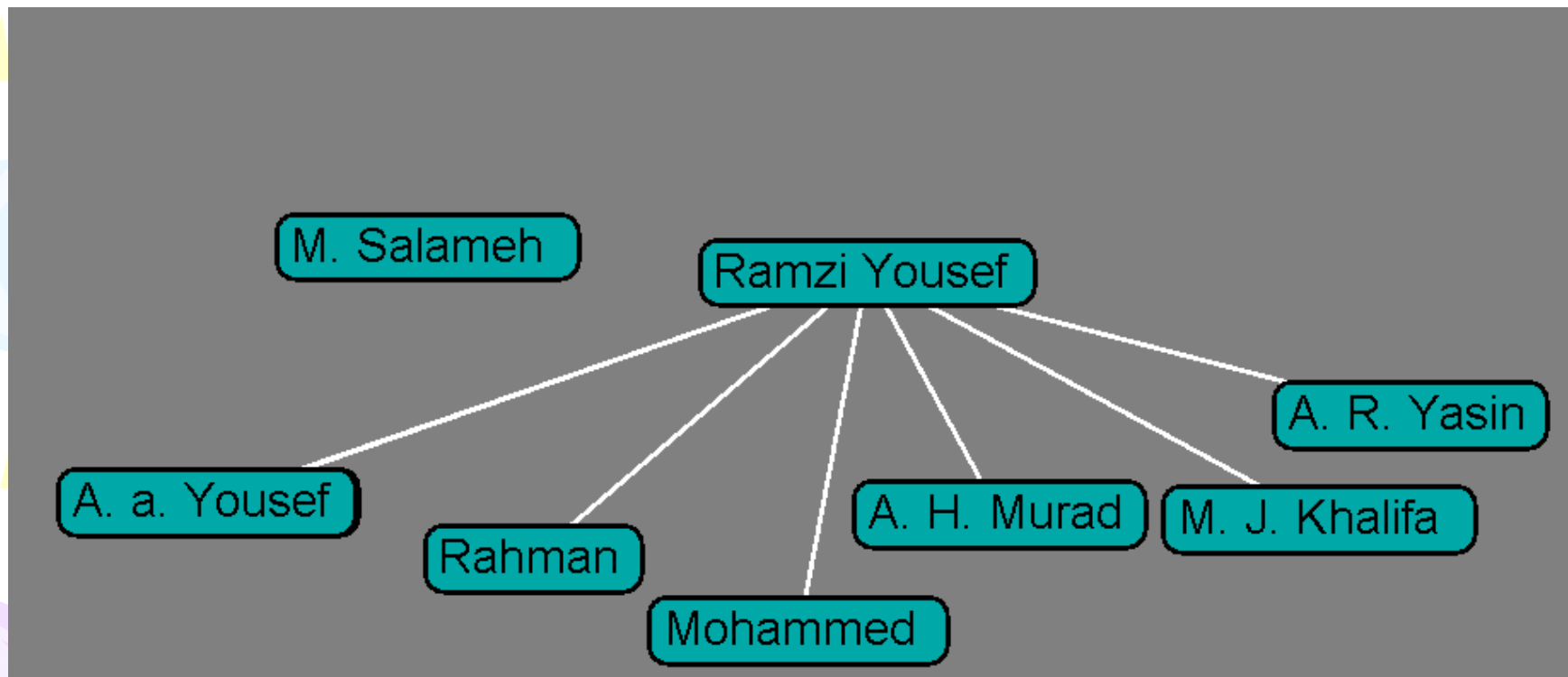
Photo By Bureau of ATF 1993 Explosives Incident Report



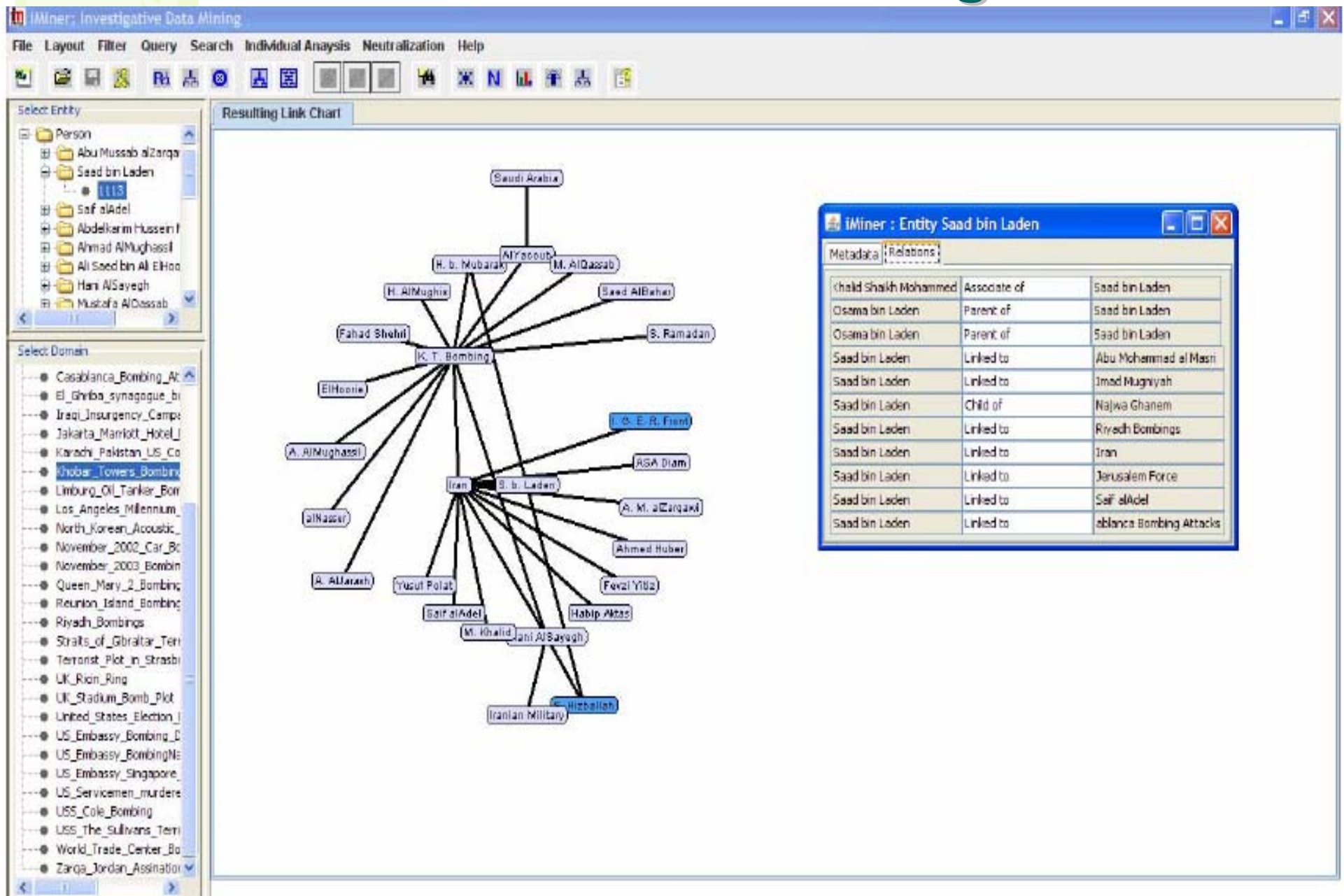
# WTC 1993 Bombing Terrorist Network



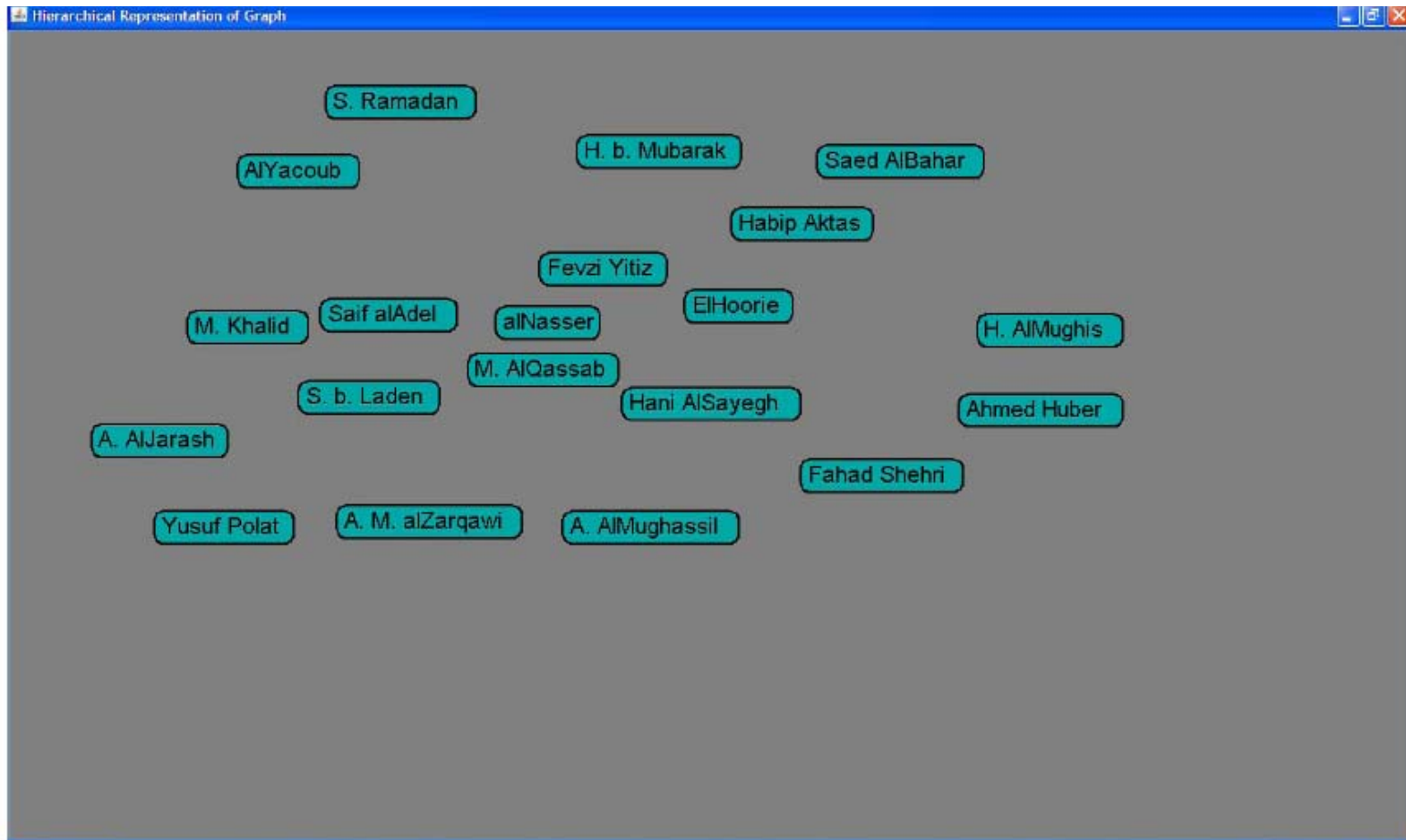
# Hierarchical Chart for WTC 1993 Bombing Terrorist Network



# Khobar Tower Bombing Plot



# No Cammand Structure found in Khabar Tower Bombing Plot





# Conclusion

- We presented an overview of Investigative Data Mining Toolkit
- The paper presented interested patterns gleaned from data. The IDM toolkit demonstrates key capabilities and concepts of a terrorist analysis toolkit.
- The toolkit can be used to understand the terrorist networks, and we are of the view that iMiner toolkit, score over traditional analysis and could reduce the consequent overload on analysts
- The model for construction of hierarchical chart illustrates the **command structure of terrorist network** and allows **destabilizing techniques** to be aimed at this command structure
- It must be noted that hierarchical chart is not a **sociogram**, but a **command structure**, used to visually demonstrate the effectiveness of targeting the highly ranked nodes to disrupt the decision making capacity
- This method attempt to **remove all command members** from the network, and leaving the **remaining cell members** without any **orders or hierarchy** in an attempt to significantly reduce the decision making capacity of the cell



Questions/ Suggestions/ comments